

ENGO 697

Remote Sensing Systems and Advanced Analytics

Session 7: How to solve inverse problems in remote sensing systems

Dr. Linlin (Lincoln) Xu

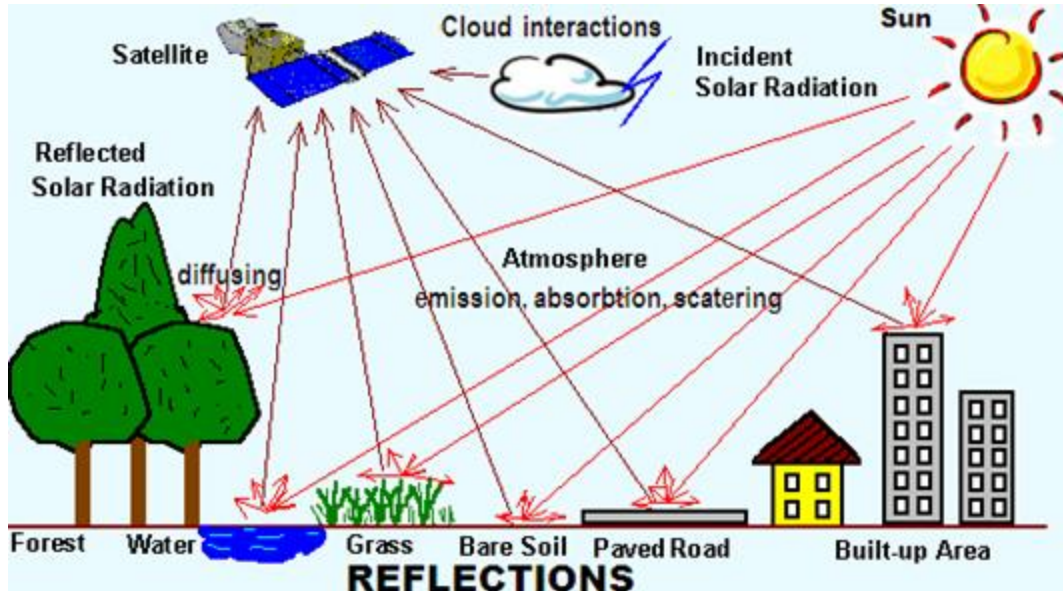
Linlin.xu@ucalgary.ca

Office: ENE 221

Outline

- Overview
- Direct inversion
- Numerical approach
- Look up table (LUT)
- Simulation & Machine learning (ML)
- Data -driven ML and DL
- Bayesian framework
- Questions

Remote Sensing System Overview



Forward model:

$$Y = f(X)$$

(1) Y: radiation received by the sensor

(2) X: variables that you want to estimate, e.g., class labels, chlorophyll content in leaves, leaf area index/density;

Inverse model:

$$X = g(Y, \theta)$$

where $g(\cdot)$ is an **unknown** inverse function with **unknown** model parameter θ .

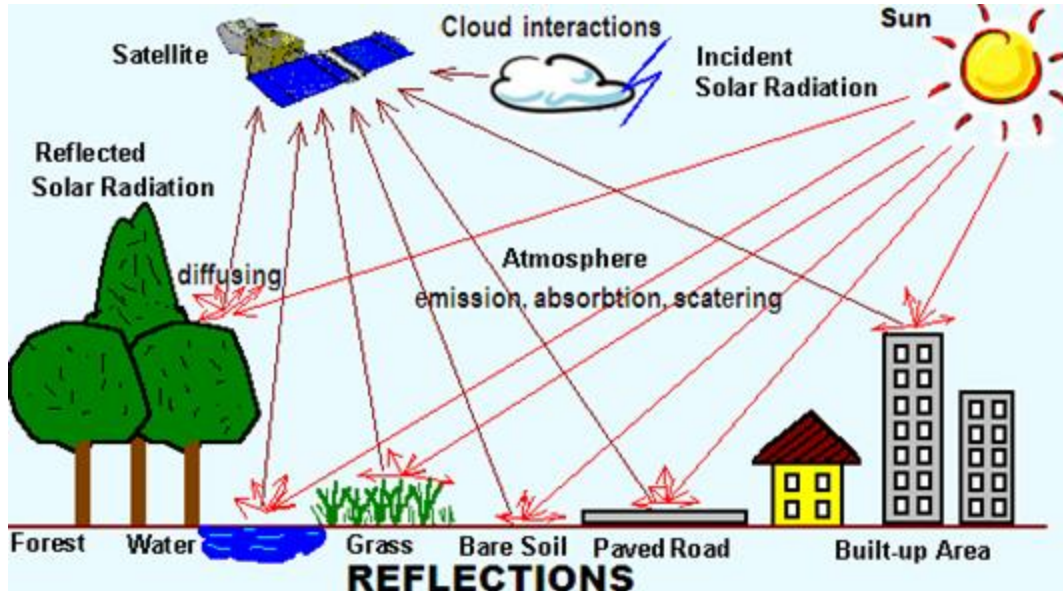
Why estimating X is difficult?

--- Knowledge $f(\cdot)$ complex, biased, highly nonlinear, with large uncertainty;

---- Data (X, Y) pairs limited, poor quality;

---- Prior information (e.g., spatial prior) ambiguous;

Remote Sensing System Overview



Forward model: $Y = f(X)$

Inverse model: $X = g(Y, \theta)$

What are the approaches for solving the inverse problems? Advantages? Disadvantages?

Common approaches for solving inverse model:

1. If $f(\cdot)$ is known and invertible, use direct inverse function approach.
1. If $f(\cdot)$ is known but highly nonlinear,
 - 2.1 use numerical approaches, e.g., Newton's approach, to estimate X using observations $\{Y\}$.
 - 2.2 use the lookup table (LUT) approach to simulate $\{(X, Y)\}$ pairs for estimating $\{X\}$ of observations $\{Y\}$.
 - 2.3 simulate $\{(X, Y)\}$ pairs to train inverse function $g(\cdot)$ which will be tested on real observations $\{Y\}$.
1. If $f(\cdot)$ is unknown, requires (Y, X) pairs, use ML and DL, to empirically build $g(\cdot)$.

(1) Direct Inverse Function Approach

In an ideal scenario, the forward model $f(\cdot)$ is simple and invertible, we can use direct inverse function approach to estimate X .

Forward model: $Y = f(X)$

where $f(\cdot)$ is invertible, $f(\cdot) \rightarrow f^{-1}(\cdot)$. For example,

Forward model: $Y = X^2$

Inverse function: $X = Y^{0.5}$

Given Y , X can be estimated based on the above inverse function.

However, in real-world remote sensing systems, the forward model $f(\cdot)$ are complicated radiative transfer models, which are usually highly nonlinear and non-invertible, and as such they do not have inverse functions.

Questions: Can this approach use data (X, Y) pairs, and spatial prior information?

Planck's Law

$$B(\lambda, T) = \frac{c_1}{\lambda^5} \left[e^{-\frac{c_2}{\lambda T}} - 1 \right]^{-1} \quad (\text{W/m}^2/\text{ster}/\mu\text{m})$$

where

λ = wavelengths in μm

T = temperature of emitting surface (deg K)

$c_1 = 1.191044 \times 10^{-5} \text{ (mW/m}^2/\text{ster/cm}^{-4}\text{)}$

$c_2 = 1.438769 \text{ (cm deg K)}$

Brightness Temperature

$$T = \frac{c_2}{\lambda \ln\left(\frac{c_1}{\lambda^5 B_\lambda} + 1\right)} \quad \text{is determined by inverting Planck function.}$$

	(1) Direct inversion	(2) LUT approach	(3) Numerical Approach	(4) Simulation & ML	(5) ML	(6) DL
$f(\cdot)$ is known	yes	yes	yes	yes	yes	yes
$f(\cdot)$ is partially known, i.e., form known, but with some unknown parameters U	no	no	Yes, estimate X and U together	no	no	no
$f(\cdot)$ unknown, (X,Y) known	no	no	no	no	yes	yes
$f(\cdot)$ unknown, (X,Y) unknown	no	no	no	no	no	no
If both $f(\cdot)$ and (X,Y) known, can accommodate both?	no	yes?	Yes? Use (X,Y) to estimate parameters in $f(\cdot)$	yes?	Yes, use both simulated and observed data	Yes, use both simulated and observed data
Can use prior information ? e.g., spatial prior and value prior	no	Yes? Use value prior for sampling	Yes? Use value prior of X in Bayesian estimation	Yes, Use value prior in sampling and spatial prior in Random fields	Yes, spatial prior in Random field approaches	Yes, similar to ML
Advantages	Knowledge-driven; Simple, easy	Knowledge-driven; Intuitive, easy, discrete fitting;	Knowledge-driven; estimate U; Efficient for simple $f(\cdot)$ in convex problems	Knowledge-driven; flexible; continuous fitting; good inter/extrapolation; faster than LUT	Data-driven; flexible; Classic;	Strong modeling capability; automatic feature learning;
Disadvantages	Unrealistic; rely on simple $f(\cdot)$	Sensitive to accuracy of $f(\cdot)$, similarity metrics, sampling density and range; slow if LUT is large; bad for extrapolation;	Rely on efficiency of nonlinear solver; Slow; Local optimum;	Overfitting and underfitting risk to simulated data; difficult model selection; Sensitive to accuracy of $f(\cdot)$, similarity metrics, sampling density and range;	Weak modeling capability; Rely on "good" engineered features; Black-box; Overfitting, underfitting; Feature and model selection is difficult and slow	Overfitting and underfitting; Black-box;

(2) Numerical Approaches

If the radiative transfer model $f(\cdot)$ is known, and we have an remote sensing observation Y , we can use the numerical approach to estimate the associated X .

Forward model: $Y = f(X)$, where $f(\cdot)$ is the radiative transfer models, which tend to be highly nonlinear and un-invertible.

Based on some observations $\{Y\}$, we can build an objective function:

$$J(X) = Y - f(X)$$

$$X = \min J(X)$$

We try to find X that through $f(\cdot)$ can generate output whose value is very close to the value of Y .

There are many methods that can solve this nonlinear optimization problems, for example,

---- Newton's method

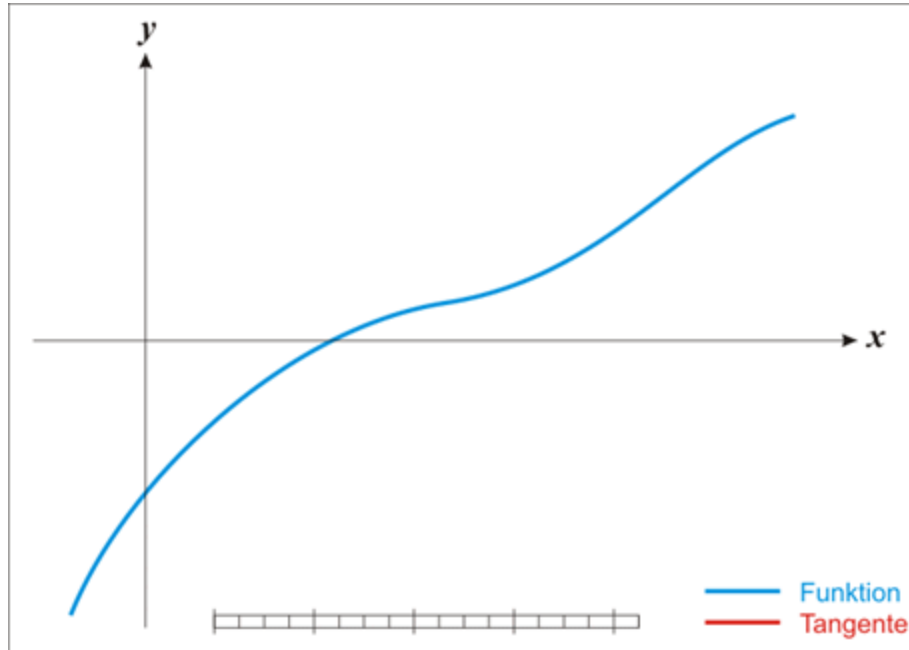
---- Gradient descent methods

---- Simulated annealing approach

Because the forward model $f(\cdot)$ contains knowledge and physical rules, this approach is usually called physical model.

Questions: How to find X that can minimize $J(X)$? Try different X and see which one gives you smallest $J(X)$? Better approach?

Root-finding Algorithm: Newton's Method



From wikipedia.org

Start with initial guess, and iteratively improve it using its tangent.

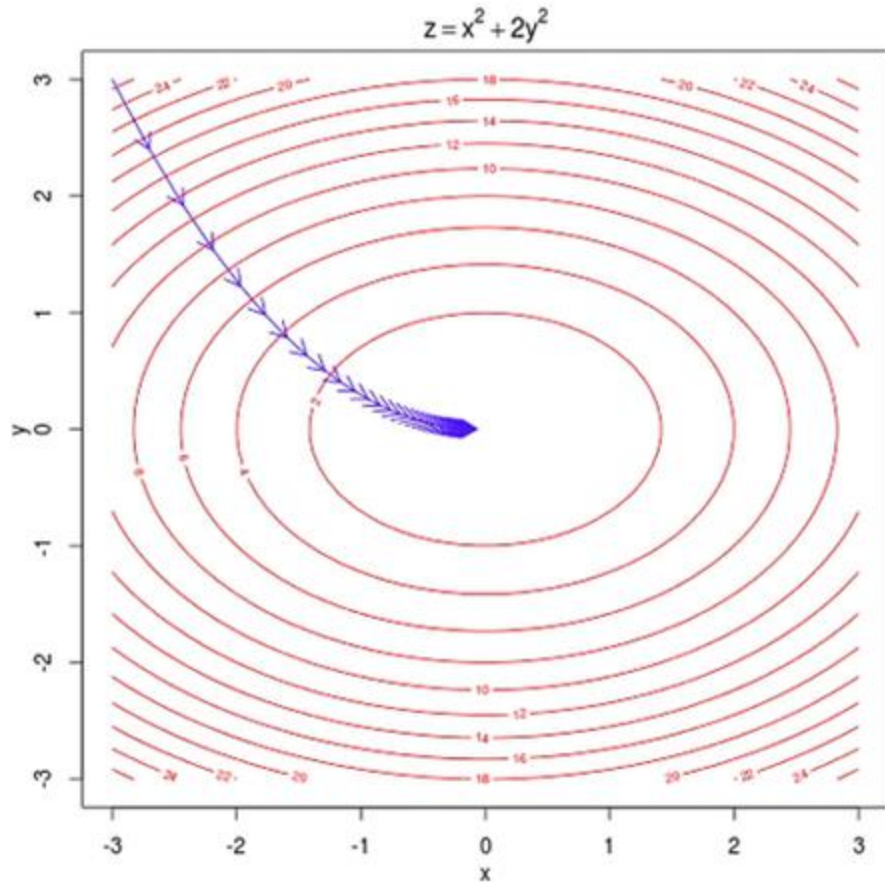
$$x_{n+1} = x_n - \frac{f(x_n)}{f'(x_n)}$$

The function f is shown in blue and the tangent line is in red.

We see that x_{n+1} is a better approximation than x_n for the root x of the function f .

The assumption is the x-intercept will typically be a better approximation to the original function's root than the previous guess, and the method can be iterated.

Gradient Descent



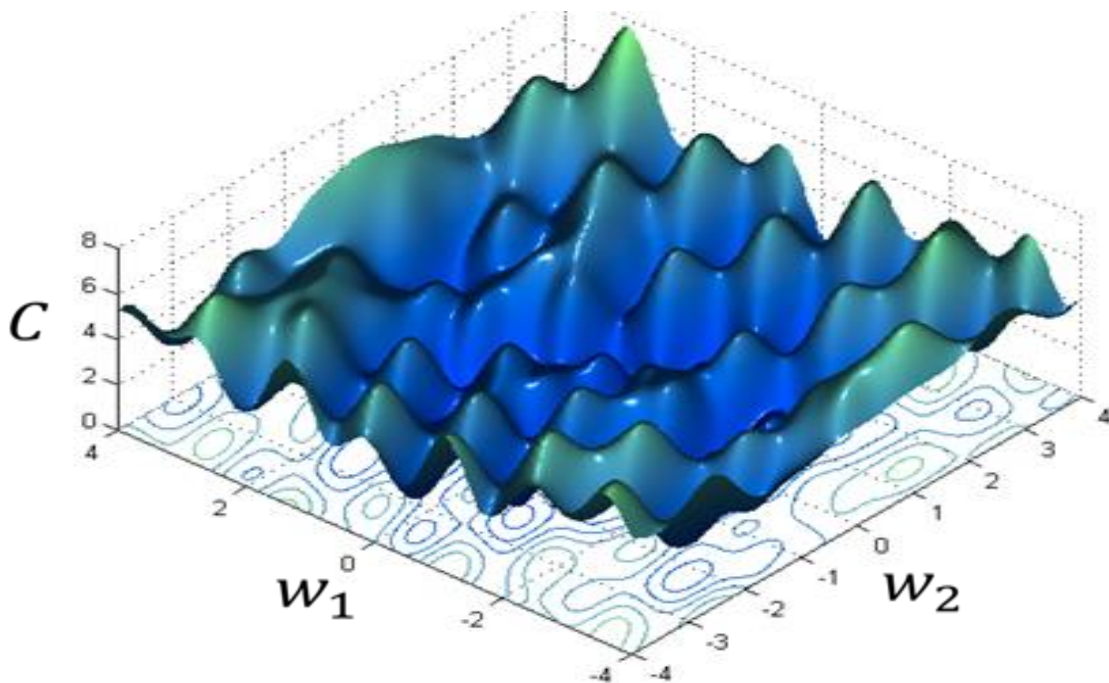
Q1: What if there are multiple local minimums for complex $f(\cdot)$?

Q2: What are roles of starting point and step length?

Start with initial guess, and iteratively improve it using its gradient.

$$\mathbf{a}_{n+1} = \mathbf{a}_n - \gamma \nabla F(\mathbf{a}_n)$$

Local Minima



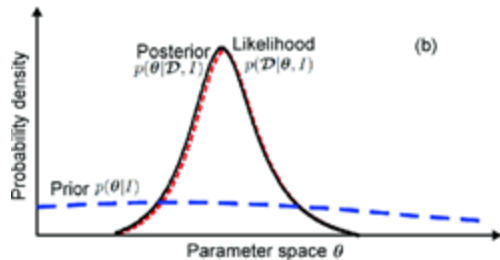
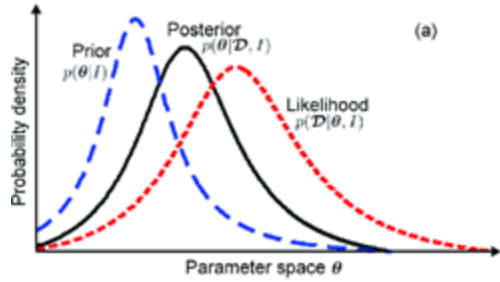
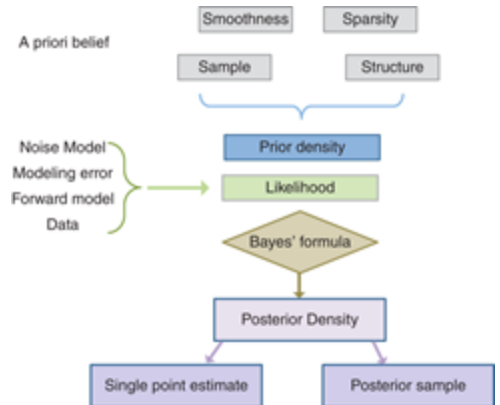
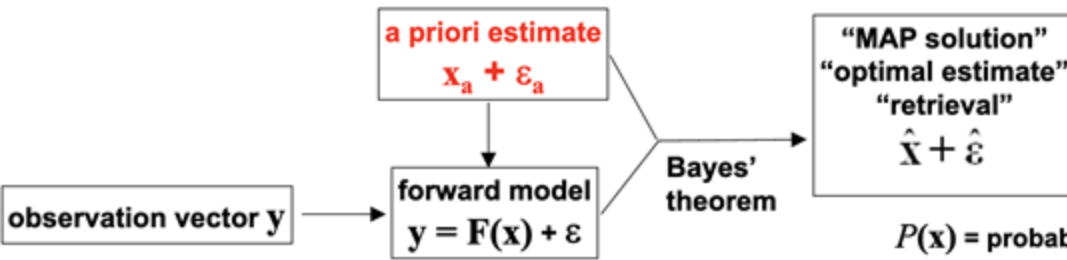
A more realistic scenario where there are many local minima;

Different initializations lead to different local minima;

How do you find global minima?

Bayesian framework - how to use "prior" information and accommodate errors?

Optimize values of an ensemble of variables (*state vector x*) using observations:



$P(x)$ = probability distribution function (pdf) of x
 $P(x,y)$ = pdf of (x,y)
 $P(y|x)$ = pdf of y given x

$$P(x,y)dxdy = P(x)dxP(y|x)dy$$

$$P(x,y)dxdy = P(y)dyP(x|y)dx$$

$$\Rightarrow P(x|y) = \frac{\overbrace{P(y|x)}^{\text{observation pdf}} \overbrace{P(x)}^{\text{a priori pdf}}}{\underbrace{P(y)}_{\text{normalizing factor (unimportant)}}}$$

Bayes' theorem

Maximum *a posteriori* (MAP) solution for x given y is defined by $\max(P(x|y))$

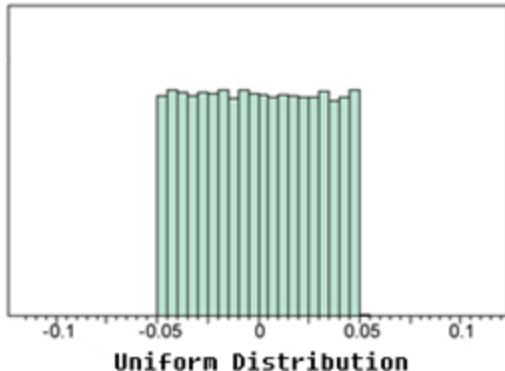
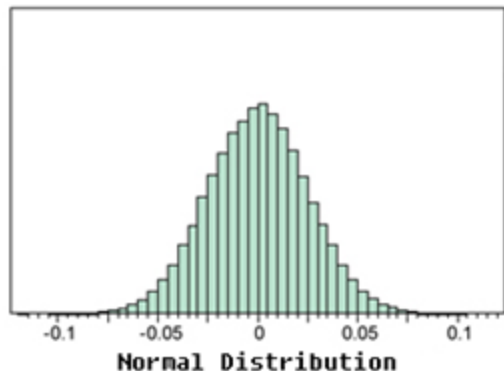
$$\Rightarrow \text{solve for } \nabla_x P(x|y) = 0$$

	(1) Direct inversion	(2) LUT approach	(3) Numerical Approach	(4) Simulation & ML	(5) ML	(6) DL
$f(\cdot)$ is known	yes	yes	yes	yes	yes	yes
$f(\cdot)$ is partially known, i.e., form known, but with some unknown parameters U	no	no	Yes, estimate X and U together	no	no	no
$f(\cdot)$ unknown, (X,Y) known	no	no	no	no	yes	yes
$f(\cdot)$ unknown, (X,Y) unknown	no	no	no	no	no	no
If both $f(\cdot)$ and (X,Y) known, can accommodate both?	no	yes?	Yes? Use (X,Y) to estimate parameters in $f(\cdot)$	yes?	Yes, use both simulated and observed data	Yes, use both simulated and observed data
Can use prior information ? e.g., spatial prior and value prior	no	Yes? Use value prior for sampling	Yes? Use value prior of X in Bayesian estimation	Yes, Use value prior in sampling and spatial prior in Random fields	Yes, spatial prior in Random field approaches	Yes, similar to ML
Advantages	Knowledge-driven; Simple, easy	Knowledge-driven; Intuitive, easy, discrete fitting;	Knowledge-driven; estimate U; Efficient for simple $f(\cdot)$ in convex problems	Knowledge-driven; flexible; continuous fitting; good inter/extrapolation; faster than LUT	Data-driven; flexible; Classic;	Strong modeling capability; automatic feature learning;
Disadvantages	Unrealistic; rely on simple $f(\cdot)$	Sensitive to accuracy of $f(\cdot)$, similarity metrics, sampling density and range; slow if LUT is large; bad for extrapolation;	Rely on efficiency of nonlinear solver; Slow; Local optimum;	Overfitting and underfitting risk to simulated data; difficult model selection; Sensitive to accuracy of $f(\cdot)$, similarity metrics, sampling density and range;	Weak modeling capability; Rely on "good" engineered features; Black-box; Overfitting, underfitting; Feature and model selection is difficult and slow	Overfitting and underfitting; Black-box;

Prosail Simulation - how to simulate data?

TABLE I: Ranges of the input variables for the PROSAIL model for the generation of the LUT.

Variable	Abbr.	Unit	Min	Max
Leaf structure parameter	N	Unitless	1	2
Leaf chlorophyll concentration	C_{ab}	μgcm^{-2}	20	70
Dry matter content	C_m	gcm^{-2}	0.004	0.007
Equivalent water thickness	C_w	gcm^{-2}	0.005	0.03
Leaf area index	LAI	m^2cm^{-2}	0.001	6
Average leaf angle	ALA	Deg	30	70
Hot-spot size parameter	hot	mm^{-1}	0.05	1
Soil brightness parameter	$scale$	Unitless	0.5	1.5



Prosail is a forward model:

$$Y = f(X)$$

What is Y?

1. Bidirectional reflectance from canopy (400nm - 2500nm);

What are the factors that constitute X:

1. A total of 14 input parameters;

How to simulate Y using X?

Step 1: know the distribution of X;

Step 2: obtain samples $\{X_i | i=1, \dots, N\}$ based on the distribution of X;

Step 3: use $\{X_i | i=1, \dots, N\}$ as input to Prosail and generate $\{Y_i | i=1, 2, \dots, N\}$

(3) LookUp Table (LUT) Approaches

If the radiative transfer model $f(\cdot)$ is known, we can simulate a collection of X and Y pairs, i.e., $\{(X_j, Y_j) \mid j=1,2,\dots,M\}$, based on which we can build a LUT and use it to estimate the X value of an observed Y value.

Forward model: $Y = f(X)$, where $f(\cdot)$ is the radiative transfer models, which tend to be highly nonlinear and not invertible.

Based on $Y=f(X)$, we build a LUT by first sampling X_j uniformly within a range $[A,B]$ (A and B are respectively the theoretical min and max value of X), and then obtaining the associated Y_j value by $Y_j=f(X_j)$.

X_1	:	Y_1
X_2	:	Y_2
X_3	:	Y_3
....		
X_n	:	Y_n

Given an observed Y value, how do we estimate its X value using LUT? First, we search the LUT to identify the row whose Y_j value is the closest to the observed Y value. Then, the X_j value associated with Y_j is treated as the estimated X value of Y .

Comparing with the numerical approaches, the LUT approach is simpler and has theoretical advantages such as being able to find the global optimum in the parameter space, and thereby the LUT approaches have been widely used in solving remote sensing inverse problems.

	(1) Direct inversion	(2) LUT approach	(3) Numerical Approach	(4) Simulation & ML	(5) ML	(6) DL
$f(\cdot)$ is known	yes	yes	yes	yes	yes	yes
$f(\cdot)$ is partially known, i.e., form known, but with some unknown parameters U	no	no	Yes, estimate X and U together	no	no	no
$f(\cdot)$ unknown, (X,Y) known	no	no	no	no	yes	yes
$f(\cdot)$ unknown, (X,Y) unknown	no	no	no	no	no	no
If both $f(\cdot)$ and (X,Y) known, can accommodate both?	no	yes?	Yes? Use (X,Y) to estimate parameters in $f(\cdot)$	yes?	Yes, use both simulated and observed data	Yes, use both simulated and observed data
Can use prior information ? e.g., spatial prior and value prior	no	Yes? Use value prior for sampling	Yes? Use value prior of X in Bayesian estimation	Yes, Use value prior in sampling and spatial prior in Random fields	Yes, spatial prior in Random field approaches	Yes, similar to ML
Advantages	Knowledge-driven; Simple, easy	Knowledge-driven; Intuitive, easy, discrete fitting;	Knowledge-driven; estimate U; Efficient for simple $f(\cdot)$ in convex problems	Knowledge-driven; flexible; continuous fitting; good inter/extrapolation; faster than LUT	Data-driven; flexible; Classic;	Strong modeling capability; automatic feature learning;
Disadvantages	Unrealistic; rely on simple $f(\cdot)$	Sensitive to accuracy of $f(\cdot)$, similarity metrics, sampling density and range; slow if LUT is large; bad for extrapolation;	Rely on efficiency of nonlinear solver; Slow; Local optimum;	Overfitting and underfitting risk to simulated data; difficult model selection; Sensitive to accuracy of $f(\cdot)$, similarity metrics, sampling density and range;	Weak modeling capability; Rely on "good" engineered features; Black-box; Overfitting, underfitting; Feature and model selection is difficult and slow	Overfitting and underfitting; Black-box;

(4) Data Simulation & Machine Learning (ML) Approaches

If the radiative transfer model $f(\cdot)$ is known, we can simulate a collection of X and Y pairs, i.e., $\{(X_j, Y_j) \mid j=1,2,\dots,M\}$. Instead of using LUT for data inversion, we can use ML approaches to learn the inverse function, i.e., $X=g(Y)$, and use this inverse function to estimate the X value of an observed Y value.

Forward model: $Y = f(X)$, where $f(\cdot)$ is the radiative transfer models, which tend to be highly nonlinear and not invertible.

Based on $Y=f(X)$, similar to LUT, we simulate the following X and Y pairs.

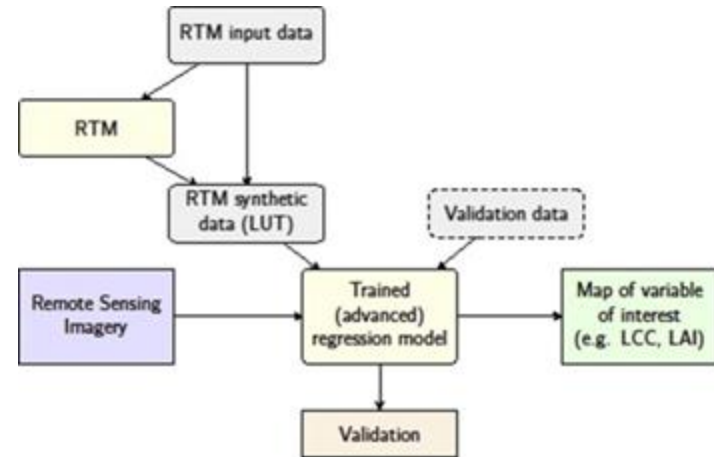
X_1 : Y_1
 X_2 : Y_2
 X_3 : Y_3
.....
 X_4 : Y_4

Based on simulated $\{(X_j, Y_j) \mid j=1,2,\dots,M\}$ pairs, we build the following objective function:

$$J(\theta) = \sum \|X_i - g(Y_i)\|$$
$$\theta = \min J(\theta)$$

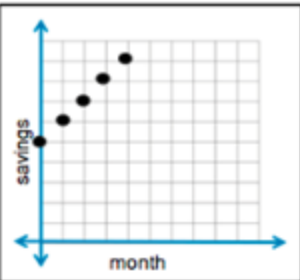
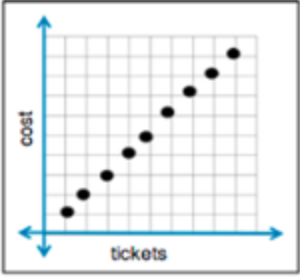
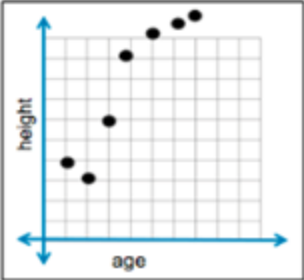
where θ is the unknown parameters in the inverse function $g(\cdot)$ which is a statistical model or machine learning model. Once we know θ , we can establish the inverse function $g(\cdot)$, and use it to estimate the X value of an observed Y value by $X=g(Y)$.

Comparing with the LUT approach that is essentially discrete interpolation, the ML approaches can learn a continuous inverse function $g(\cdot)$ using the simulated data, and thereby they theoretically can achieve more accurate estimation. Moreover, ML approaches tend to be faster because they do not need to do LUT searching for every observation.

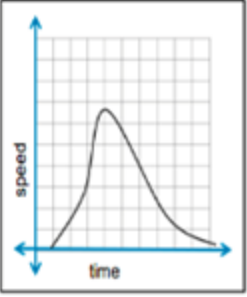
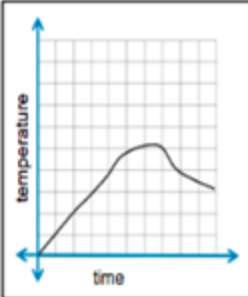
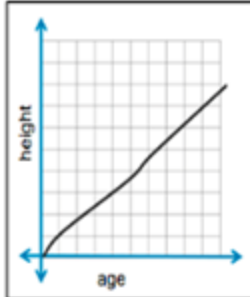
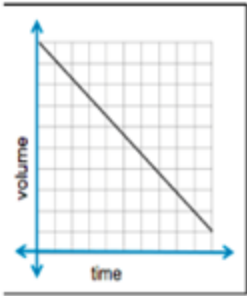
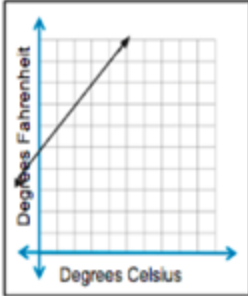
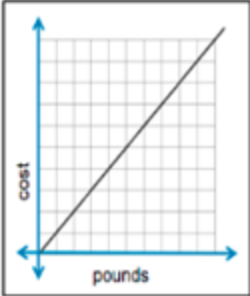


LUT vs. Machine learning

DISCRETE



CONTINUOUS



True inverse function vs. approximated inverse function

Forward model:

$$Y = f(X)$$

(1) Y: received radiation by the sensor

(2) X: variables that you want to know, e.g., class labels, chlorophyll content in leaves, leaf area index/density;

True inverse function:

$$X = t(Y) = f^{-1}(Y)$$

where $f^{-1}(\cdot)$ is difficult to get and the form of $t(\cdot)$ is usually unknown;

Approximated inverse function:

$$X = g(Y)$$

Note that $g(\cdot)$ is only an approximation to the true inverse function $t(\cdot)$

Appropriate fitting: when the complexity of $g(\cdot)$ is close to $t(\cdot)$;

Overfitting: when the complexity of $g(\cdot)$ is larger than $t(\cdot)$;

Underfitting: when the complexity of $g(\cdot)$ is smaller than $t(\cdot)$;



Based on $\{(X_j, Y_j) \mid j=1,2,\dots,n\}$, we build the following objective function:

$$J(\theta) = \sum \|X_i - g(Y_i)\|$$

$$\theta = \min J(\theta)$$

Planck's Law

where

$$B(\lambda, T) = \frac{c_1}{\lambda^5} \left[e^{-\frac{c_2}{\lambda T}} - 1 \right]^{-1} \quad (\text{W/m}^2/\text{ster}/\mu\text{m})$$

λ = wavelengths in μm

T = temperature of emitting surface (deg K)

$c_1 = 1.191044 \times 10^{-5}$ (mW/m²/ster/cm⁻⁴)

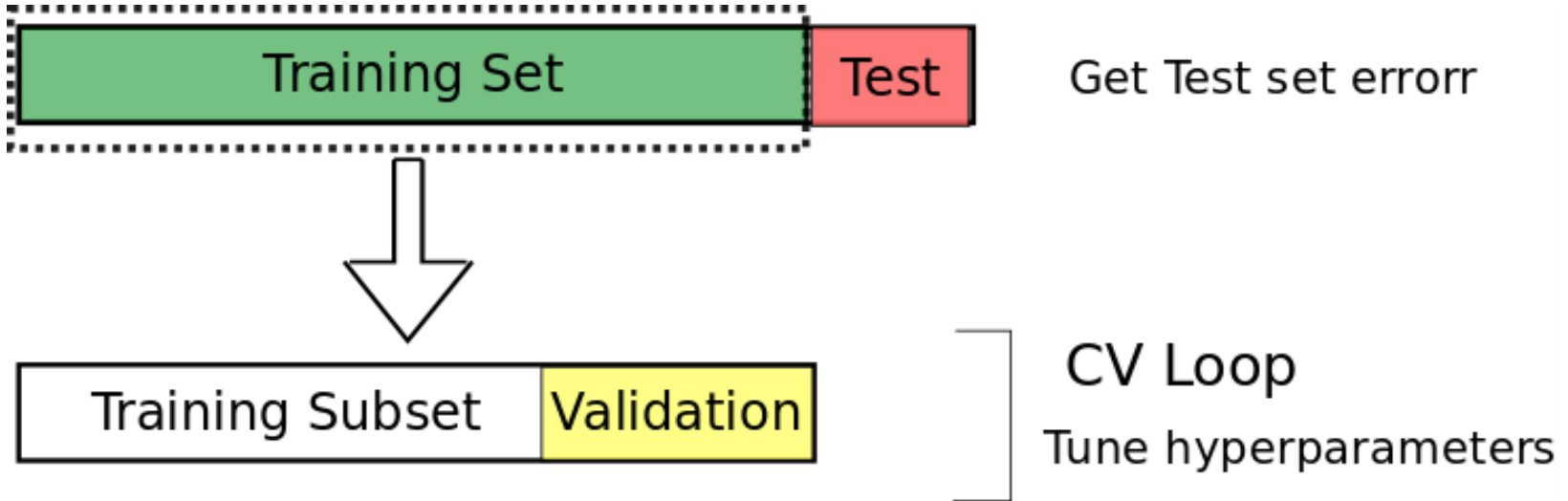
$c_2 = 1.438769$ (cm deg K)

What happens if you fit the BT function using a million-layers neural network?

Brightness Temperature

$$T = \frac{c_1}{\lambda^5 B_\lambda} \left[\lambda \ln \left(\frac{c_1}{\lambda^5 B_\lambda} + 1 \right) \right] \text{ is determined by inverting Planck function.}$$

How do you select the “best” $g(\cdot)$?



Try different models, $g_1(\cdot)$, $g_2(\cdot)$, ..., $g_n(\cdot)$, and select the one that with highest accuracy on **the validation set**.

Overfitting vs. Underfitting

Overfitting:

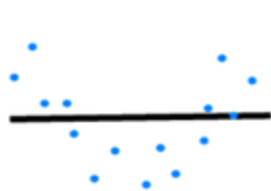
---- ML model is so **flexible and complex** that it **accommodates the noise effect** in the training data and treats it as signal, and **the learnt noise characteristics** cannot generalize well to the test data;

---- **very high training accuracy but low validation/test accuracy;**

Underfitting:

---- ML model is so **simple and rigid** that it **does not has enough capacity to accommodate signal** in the training data, and the **learnt biased/partial information** cannot generalize well to the test data;

---- **low training accuracy & low validation/test accuracy;**



Underfitting

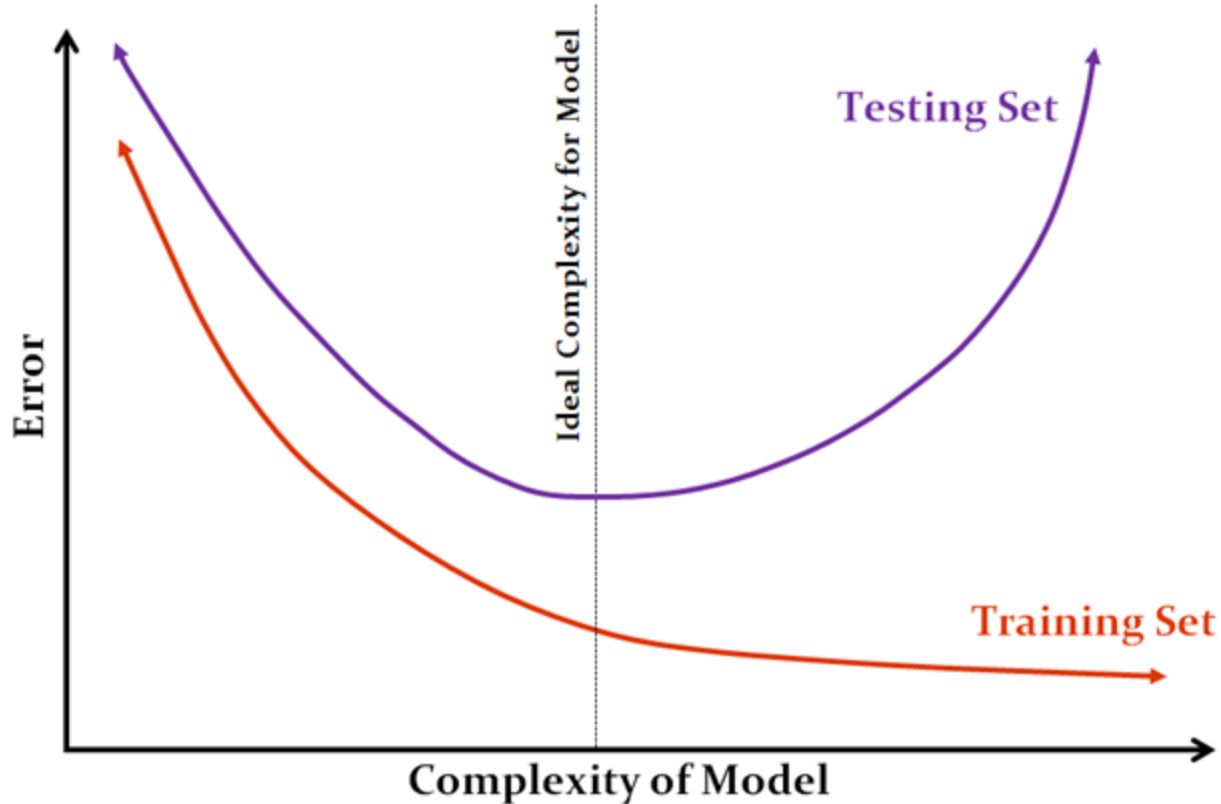


Desired

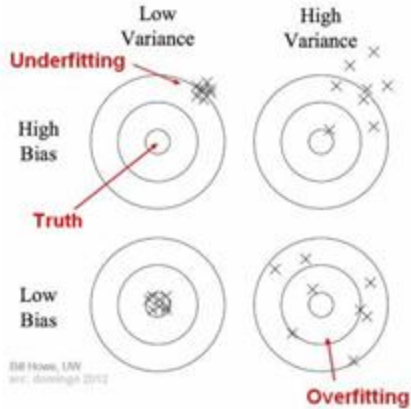


Overfitting

Training error vs. test error as model complexity changes



What are Bias and Variance? Why overfitting means small bias and big variance?



Bias is the **difference** between the **average prediction** of our model and the **true value** which we are trying to predict.

Variance is the **variability** of model **prediction**.

Why increasing model complexity lead to small bias in prediction?

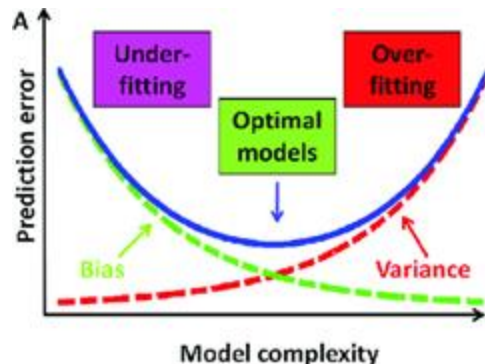
---- **increasing model complexity** -> $g(\cdot)$ to be **universal approximator** -> stronger **accommodating/modeling capability** to learn the **genuine nonlinear relationship between X and Y** in $X = t(Y)$ -> **less bias**;

Why increasing model complexity lead to larger variance in prediction?

---- **increasing model complexity** -> $g(\cdot)$ to be **universal approximator** -> stronger **accommodating/modeling capability** to learn both the **genuine nonlinear relationship between X and Y** and **irrelevant factors (i.e., noise and even errors in the data)** -> **larger variance**;

Why decreasing model complexity lead to larger bias in prediction?

Why increasing model complexity lead to smaller variance in prediction?



Overfitting vs. Underfitting

Overfitting:

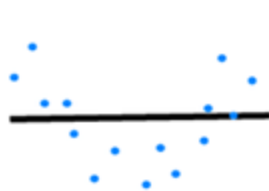
---- ML model is so **flexible and complex** that it **accommodates the noise effect** in the training data and treats it as signal, and **the learnt noise characteristics** cannot generalize well to the test data;

---- **very high training accuracy but low validation/test accuracy**; **small Bias but big variation** in prediction;

Underfitting:

---- ML model is so **simple and rigid** that it **does not has enough capacity to accommodate signal** in the training data, and the **learnt biased/parcial information** cannot generalize well to the test data;

---- **low training accuracy & low validation/test accuracy**; **big Bias but small variation** in prediction;



Underfitting



Desired



Overfitting

	(1) Direct inversion	(2) LUT approach	(3) Numerical Approach	(4) Simulation & ML	(5) ML	(6) DL
$f(\cdot)$ is known	yes	yes	yes	yes	yes	yes
$f(\cdot)$ is partially known, i.e., form known, but with some unknown parameters U	no	no	Yes, estimate X and U together	no	no	no
$f(\cdot)$ unknown, (X,Y) known	no	no	no	no	yes	yes
$f(\cdot)$ unknown, (X,Y) unknown	no	no	no	no	no	no
If both $f(\cdot)$ and (X,Y) known, can accommodate both?	no	yes?	Yes? Use (X,Y) to estimate parameters in $f(\cdot)$	yes?	Yes, use both simulated and observed data	Yes, use both simulated and observed data
Can use prior information ? e.g., spatial prior and value prior	no	Yes? Use value prior for sampling	Yes? Use value prior of X in Bayesian estimation	Yes, Use value prior in sampling and spatial prior in Random fields	Yes, spatial prior in Random field approaches	Yes, similar to ML
Advantages	Knowledge-driven; Simple, easy	Knowledge-driven; Intuitive, easy, discrete fitting;	Knowledge-driven; estimate U; Efficient for simple $f(\cdot)$ in convex problems	Knowledge-driven; flexible; continuous fitting; good inter/extrapolation; faster than LUT	Data-driven; flexible; Classic;	Strong modeling capability; automatic feature learning;
Disadvantages	Unrealistic; rely on simple $f(\cdot)$	Sensitive to accuracy of $f(\cdot)$, similarity metrics, sampling density and range; slow if LUT is large; bad for extrapolation;	Rely on efficiency of nonlinear solver; Slow; Local optimum;	Overfitting and underfitting risk to simulated data; difficult model selection; Sensitive to accuracy of $f(\cdot)$, similarity metrics, sampling density and range;	Weak modeling capability; Rely on "good" engineered features; Black-box; Overfitting, underfitting; Feature and model selection is difficult and slow	Overfitting and underfitting; Black-box;

(5) Machine Learning (ML) Approaches

If $f(\cdot)$ is unknown, we resort to data for solving inverse problems. We need to collect X and Y measurements, i.e., $\{(X_j, Y_j) \mid j=1,2,\dots,T\}$, based on which we establish the inverse function $X=g(Y, \theta)$, where $g(\cdot)$ is a statistical or ML model (empirical models).

No $f(X)$

We need to measure Y and the associated ground truth data X to get some (X, Y) pairs, to train ML model.

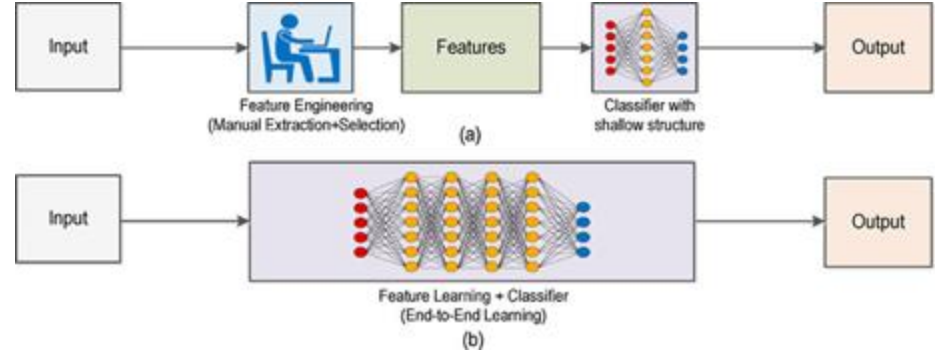
$X_1 : Y_1$

$X_2 : Y_2$

$X_3 : Y_3$

....

$X_4 : Y_4$



Based on $\{(X_j, Y_j) \mid j=1,2,\dots,T\}$, we build the following objective function:

$$J(\theta) = \sum \|X_i - g(Y_i)\|$$

$$\theta = \min J(\theta)$$

where θ is the unknown parameters in $g(\cdot)$. Once we know θ , we can establish the inverse function $g(\cdot)$, and use it to estimate the X value of an observed Y value by $X=g(Y)$.

Limitations of ML approaches, when X and Y have highly-nonlinear relationship? Why feature extraction and selection? Why not using original features?

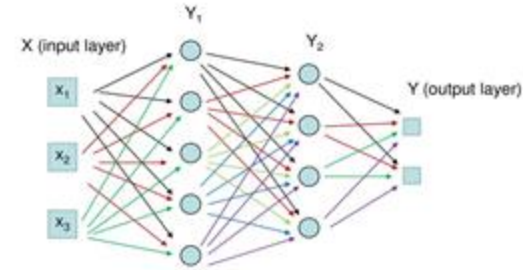
	(1) Direct inversion	(2) LUT approach	(3) Numerical Approach	(4) Simulation & ML	(5) ML	(6) DL
$f(\cdot)$ is known	yes	yes	yes	yes	yes	yes
$f(\cdot)$ is partially known, i.e., form known, but with some unknown parameters U	no	no	Yes, estimate X and U together	no	no	no
$f(\cdot)$ unknown, (X,Y) known	no	no	no	no	yes	yes
$f(\cdot)$ unknown, (X,Y) unknown	no	no	no	no	no	no
If both $f(\cdot)$ and (X,Y) known, can accommodate both?	no	yes?	Yes? Use (X,Y) to estimate parameters in $f(\cdot)$	yes?	Yes, use both simulated and observed data	Yes, use both simulated and observed data
Can use prior information ? e.g., spatial prior and value prior	no	Yes? Use value prior for sampling	Yes? Use value prior of X in Bayesian estimation	Yes, Use value prior in sampling and spatial prior in Random fields	Yes, spatial prior in Random field approaches	Yes, similar to ML
Advantages	Knowledge-driven; Simple, easy	Knowledge-driven; Intuitive, easy, discrete fitting;	Knowledge-driven; estimate U; Efficient for simple $f(\cdot)$ in convex problems	Knowledge-driven; flexible; continuous fitting; good inter/extrapolation; faster than LUT	Data-driven; flexible; Classic;	Strong modeling capability; automatic feature learning;
Disadvantages	Unrealistic; rely on simple $f(\cdot)$	Sensitive to accuracy of $f(\cdot)$, similarity metrics, sampling density and range; slow if LUT is large; bad for extrapolation;	Rely on efficiency of nonlinear solver; Slow; Local optimum;	Overfitting and underfitting risk to simulated data; difficult model selection; Sensitive to accuracy of $f(\cdot)$, similarity metrics, sampling density and range;	Weak modeling capability; Rely on "good" engineered features; Black-box; Overfitting, underfitting; Feature and model selection is difficult and slow	Overfitting and underfitting; Black-box;

(6) Deep Learning (DL) Approaches

Deep learning (DL) approaches are also ML approaches, and as such they can be used for data inversion through (4) and (5), i.e.,

- if $f(\cdot)$ is known, we simulate $\{(X_j, Y_j) \mid j=1,2,\dots,T\}$ using $f(\cdot)$ and use them to train DL models for obtaining the inverse function $X=g(Y)$;
- if $f(\cdot)$ is unknown, we obtain remote sensing data Y and ground truth data X to build X and Y pairs, i.e., $\{(X_j, Y_j) \mid j=1,2,\dots,T\}$, and use them to train DL models for obtaining the inverse function $X=g(Y)$;

(feedforward) Neural network



Based on **training data**, we build the following objective function:

$$J(\theta) = \sum \|X_i - g(Y_i)\|$$
$$\theta = \min J(\theta)$$

$$Y_1 = F(W^{(1)}X + b^{(1)}), \quad Y_2 = F(W^{(2)}Y_1 + b^{(2)}), \quad Y = W^{(3)}Y_2 + b^{(3)}$$

where θ is the unknown parameters in DL model $g(\cdot)$. Once we know θ , we can establish the inverse function $g(\cdot)$, and use it to estimate the X value of an observed Y value by $X=g(Y)$.

Comparing with traditional ML approaches, such as SVM and random forest, the DL approaches, due to their strong modeling capability and GPU computation, are more capable of effectively and efficiently learning the complex nonlinear relationship between Y and X , and perform accurate and fast model prediction for estimating X .

When DL is better than ML?

Forward model:

$$Y = f(X)$$

(1) Y: received radiation by the sensor

(2) X: variables that you want to know, e.g., class labels, chlorophyll content in leaves, leaf area index/density;

True inverse function:

$$X = t(Y) = f^{-1}(Y)$$

where $f^{-1}(\cdot)$ is difficult to get and the form of $t(\cdot)$ is usually unknown;

Approximated inverse function:

$$X = g(Y)$$

Note that $g(\cdot)$ is only an approximation to the true inverse function $t(\cdot)$.

Appropriate fitting: when the complexity of $g(\cdot)$ is close to $t(\cdot)$;

Overfitting: when the complexity of $g(\cdot)$ is larger than $t(\cdot)$;

Underfitting: when the complexity of $g(\cdot)$ is smaller than $t(\cdot)$;

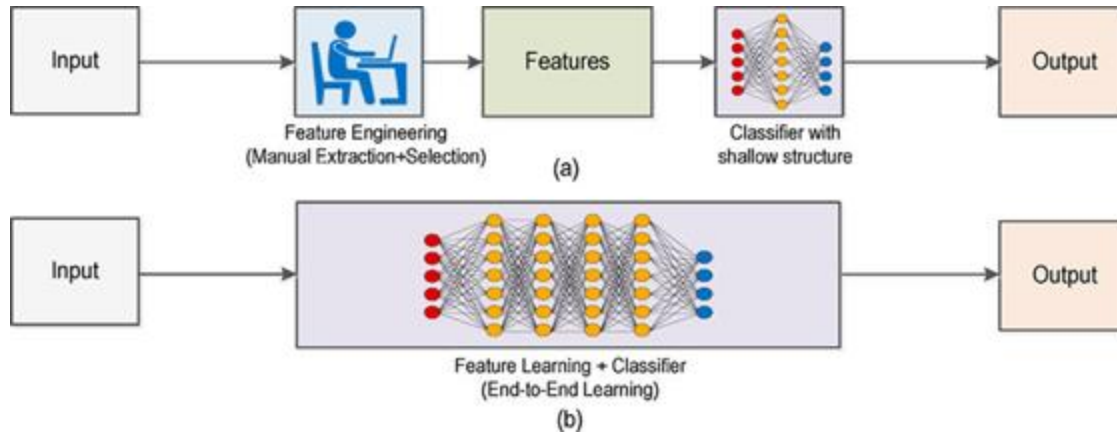


Based on $\{(X_j, Y_j) \mid j=1,2,\dots,n\}$, we build the following objective function:

$$J(\theta) = \sum \|X_i - g(Y_i)\|$$

$$\theta = \min J(\theta)$$

Feature-driven machine learning vs. Data-driven deep learning (DL)



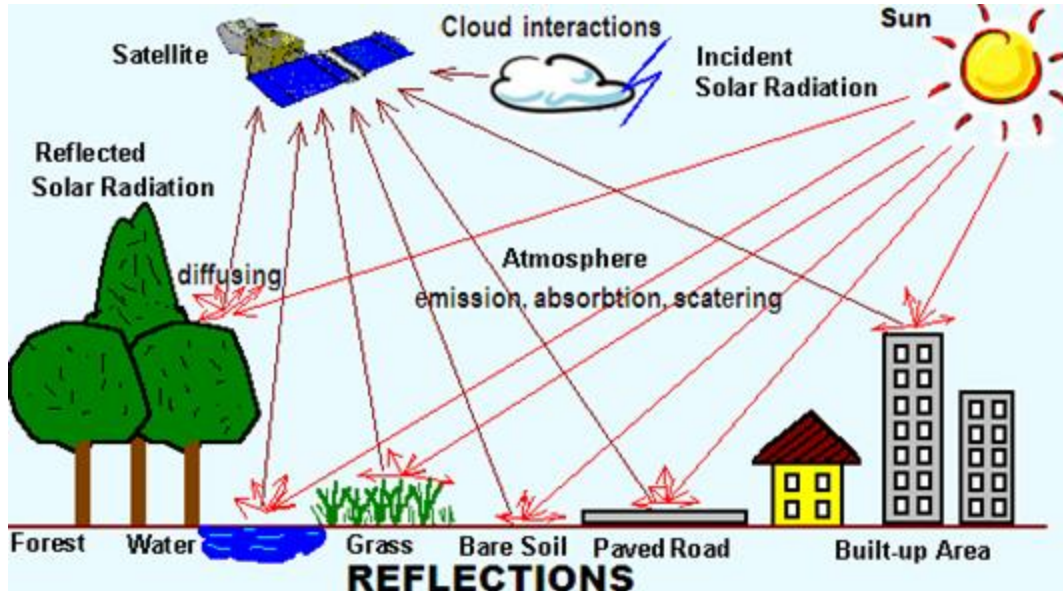
Advantages of DL approaches for RS image classification:

- (1) automatically learn the “best” feature without requiring task-specific classifier-specific knowledge;
- (2) End-to-end approach without any intermediate stages in the data-processing pipeline;
- (3) Complex model -> strong modeling capability -> efficiently capture the subtle differences among classes;
- (4) Powerful GPU computation

	(1) Direct inversion	(2) LUT approach	(3) Numerical Approach	(4) Simulation & ML	(5) ML	(6) DL
$f(\cdot)$ is known	yes	yes	yes	yes	yes	yes
$f(\cdot)$ is partially known, i.e., form known, but with some unknown parameters U	no	no	Yes, estimate X and U together	no	no	no
$f(\cdot)$ unknown, (X,Y) known	no	no	no	no	yes	yes
$f(\cdot)$ unknown, (X,Y) unknown	no	no	no	no	no	no
If both $f(\cdot)$ and (X,Y) known, can accommodate both?	no	yes?	Yes? Use (X,Y) to estimate parameters in $f(\cdot)$	yes?	Yes, use both simulated and observed data	Yes, use both simulated and observed data
Can use prior information ? e.g., spatial prior and value prior	no	Yes? Use value prior for sampling	Yes? Use value prior of X in Bayesian estimation	Yes, Use value prior in sampling and spatial prior in Random fields	Yes, spatial prior in Random field approaches	Yes, similar to ML
Advantages	Knowledge-driven; Simple, easy	Knowledge-driven; Intuitive, easy, discrete fitting;	Knowledge-driven; estimate U; Efficient for simple $f(\cdot)$ in convex problems	Knowledge-driven; flexible; continuous fitting; good inter/extrapolation; faster than LUT	Data-driven; flexible; Classic;	Strong modeling capability; automatic feature learning;
Disadvantages	Unrealistic; rely on simple $f(\cdot)$	Sensitive to accuracy of $f(\cdot)$, similarity metrics, sampling density and range; slow if LUT is large; bad for extrapolation;	Rely on efficiency of nonlinear solver; Slow; Local optimum;	Overfitting and underfitting risk to simulated data; difficult model selection; Sensitive to accuracy of $f(\cdot)$, similarity metrics, sampling density and range;	Weak modeling capability; Rely on "good" engineered features; Black-box; Overfitting, underfitting; Feature and model selection is difficult and slow	Overfitting and underfitting; Black-box;

More on Error and Prior

Remote Sensing System Overview



Forward model:

$$Y = f(X)$$

(1) Y: radiation received by the sensor

(2) X: variables that you want to estimate, e.g., class labels, chlorophyll content in leaves, leaf area index/density;

Inverse model:

$$X = g(Y, \theta)$$

where $g(\cdot)$ is an **unknown** inverse function with **unknown** model parameter θ .

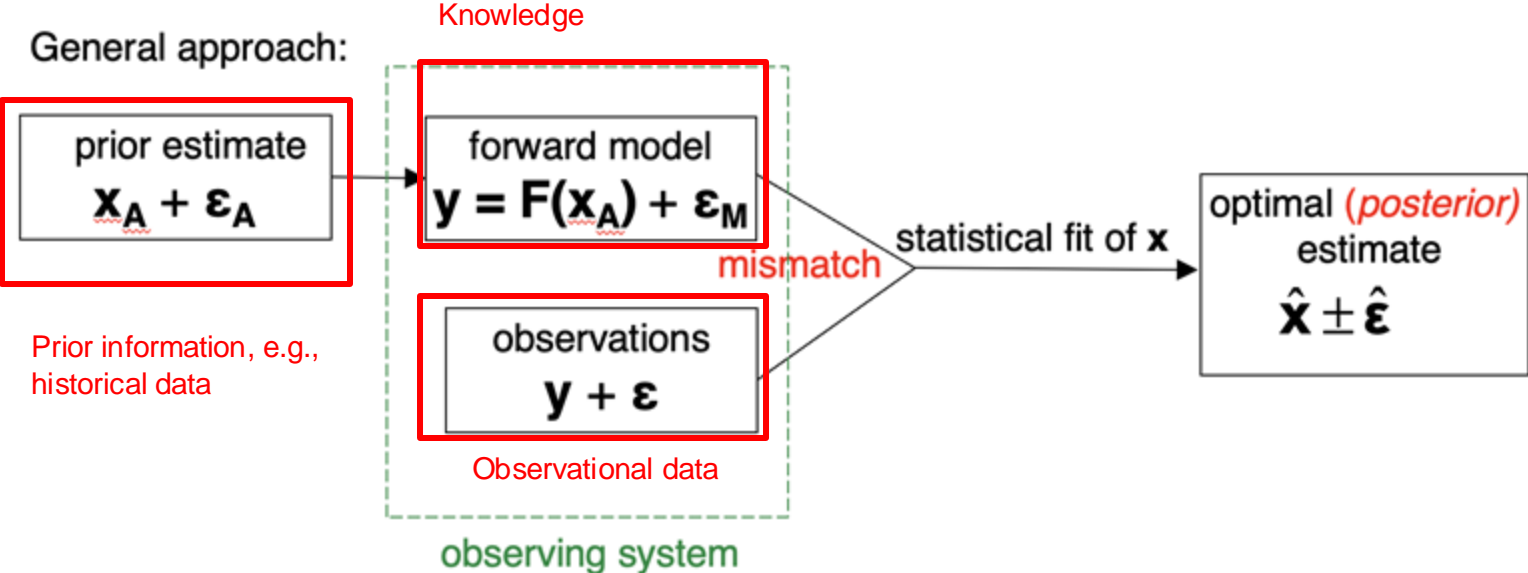
Why estimating X is difficult?

- Knowledge $f(\cdot)$ complex, biased, highly nonlinear, with large uncertainty;
- Data (X, Y) pairs limited, poor quality;
- Prior information (e.g., spatial prior) ambiguous;

How to use prior estimate (e.g., historical estimate) and accommodate errors?

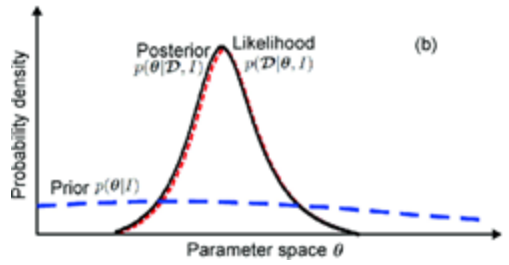
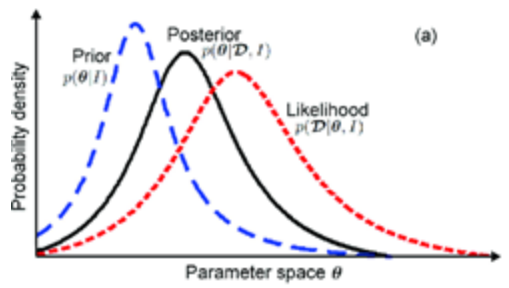
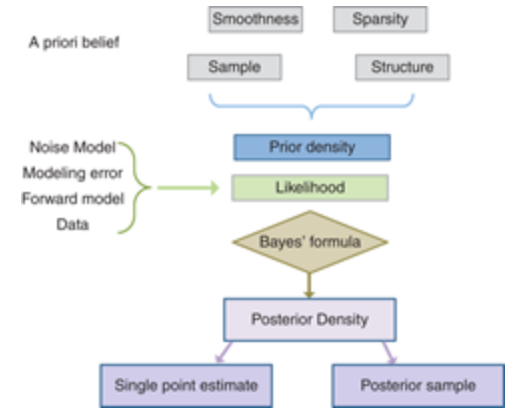
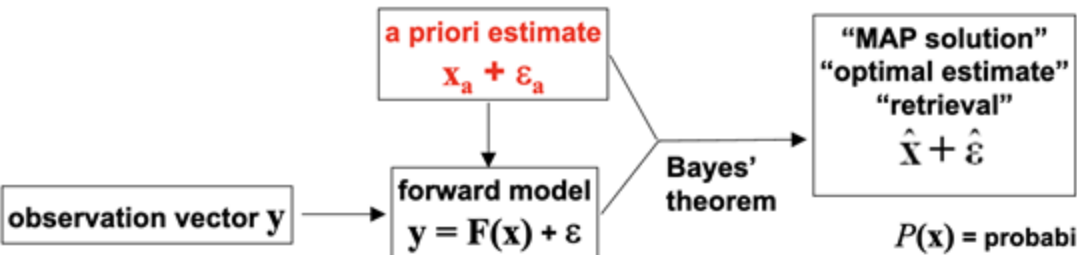
Estimate selected variables driving a physical system (*state vector \mathbf{x} , dim n*) by using:

- the observable manifestations of the system (*observation vector \mathbf{y} , dim m*)
- a physical model $\mathbf{y} = \mathbf{F}(\mathbf{x})$ (*forward model*)
- a *prior estimate \mathbf{x}_A* before the observations have been made
- Statistics for the errors $\boldsymbol{\varepsilon}$ in the different components of the problem



Bayesian framework - use “prior” & conditional distributions

Optimize values of an ensemble of variables (*state vector x*) using observations:



$P(x)$ = probability distribution function (pdf) of x
 $P(x,y)$ = pdf of (x,y)
 $P(y|x)$ = pdf of y given x

$$P(x,y)dxdy = P(x)dxP(y|x)dy$$

$$P(x,y)dxdy = P(y)dyP(x|y)dx$$

$$\Rightarrow \underbrace{P(x|y)}_{\text{a posteriori pdf}} = \frac{\underbrace{P(y|x)}_{\text{observation pdf}} \underbrace{P(x)}_{\text{a priori pdf}}}{\underbrace{P(y)}_{\text{normalizing factor (unimportant)}}}$$

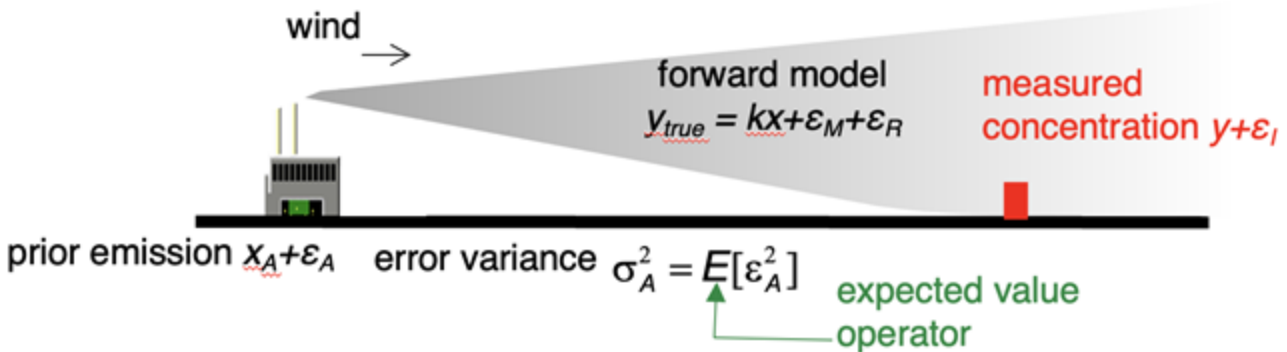
Bayes' theorem

Maximum a posteriori (MAP) solution for x given y is defined by $\max(P(x|y))$

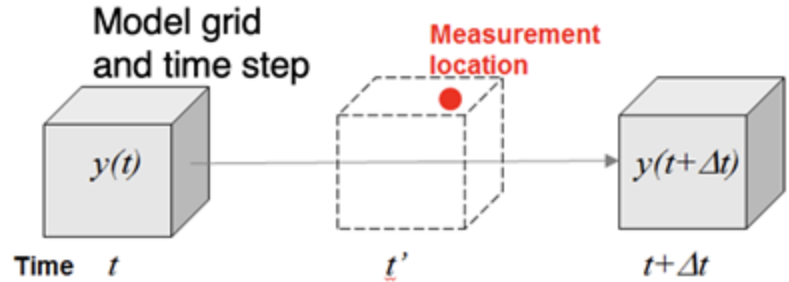
$$\Rightarrow \text{solve for } \nabla_x P(x|y) = 0$$

A simple linear inverse problem

Optimize estimate of true emission x using a single concentration measurement y ; assume simple linear transport forward model described by $y = kx$



Instrument error ϵ_I
 Forward model error ϵ_M
 Representation error ϵ_R :



Observational error $\epsilon_O = \epsilon_I + \epsilon_M + \epsilon_R$ $y = kx + \epsilon_O$ error variance $\sigma_O^2 = E[\epsilon_O^2]$

Error variances are additive: $\sigma_O^2 = \sigma_I^2 + \sigma_M^2 + \sigma_R^2$

How to solve this inverse problem using Bayesian MAP estimation? error distribution assumption?

Prior pdf:
$$P(x) = \frac{1}{\sigma_A \sqrt{2\pi}} \exp\left[-\frac{(x - x_A)^2}{2\sigma_A^2}\right]$$

Observation pdf:
$$P(y|x) = \frac{1}{\sigma_o \sqrt{2\pi}} \exp\left[-\frac{(y - kx)^2}{2\sigma_o^2}\right]$$

Posterior pdf:
$$P(x|y) \propto P(x)P(y|x) \propto \exp\left[-\frac{(x - x_A)^2}{2\sigma_A^2} - \frac{\boxed{}}{2\sigma_o^2}\right]$$

Cost function to minimize:
$$J(x) = \frac{(x - x_A)^2}{\sigma_A^2} + \frac{(y - kx)^2}{\sigma_o^2} \quad \text{solve } \underline{dJ/dx} = 0$$

Posterior solution:
$$\hat{x} = x_A + g(y - kx_A) \quad g = \frac{k\sigma_A^2}{k^2\sigma_A^2 + \sigma_o^2} \quad \text{gain factor}$$

with posterior pdf
$$P(x|y) = \frac{1}{\hat{\sigma} \sqrt{2\pi}} \exp\left[-\frac{(x - \hat{x})^2}{2\hat{\sigma}^2}\right] \quad \frac{1}{\hat{\sigma}^2} = \frac{1}{\sigma_A^2} + \frac{1}{(\sigma_o/k)^2}$$

posterior error variance

Relationship to truth:
$$\hat{x} = \underset{\text{optimal estimate}}{x} + (1 - a) \underset{\text{truth}}{(x_A - x)} + \underset{\text{smoothing error}}{g} \underset{\text{obs. error}}{\varepsilon_o}$$

 $a = gk = \frac{\sigma_A^2}{\sigma_A^2 + (\sigma_o/k)^2}$
averaging kernel

How to solve this inverse problem using Bayesian MAP estimation?

Prior pdf:
$$P(x) = \frac{1}{\sigma_A \sqrt{2\pi}} \exp\left[-\frac{(x - x_A)^2}{2\sigma_A^2}\right]$$

Observation pdf:
$$P(y|x) = \frac{1}{\sigma_o \sqrt{2\pi}} \exp\left[-\frac{(y - kx)^2}{2\sigma_o^2}\right]$$

Posterior pdf:
$$P(x|y) \propto P(x)P(y|x) \propto \exp\left[-\frac{(x - x_A)^2}{2\sigma_A^2} - \frac{(y - kx)^2}{2\sigma_o^2}\right]$$

Cost function to minimize:
$$J(x) = \frac{(x - x_A)^2}{\sigma_A^2} + \frac{(y - kx)^2}{\sigma_o^2} \quad \text{solve } \underline{dJ/dx} = 0$$

Posterior solution:
$$\hat{x} = x_A + g(y - kx_A) \quad g = \frac{k\sigma_A^2}{k^2\sigma_A^2 + \sigma_o^2} \quad \text{gain factor}$$

with posterior pdf
$$P(x|y) = \frac{1}{\hat{\sigma} \sqrt{2\pi}} \exp\left[-\frac{(x - \hat{x})^2}{2\hat{\sigma}^2}\right] \quad \frac{1}{\hat{\sigma}^2} = \frac{1}{\sigma_A^2} + \frac{1}{(\sigma_o/k)^2}$$

posterior error variance

Relationship to truth:
$$\hat{x} = \underbrace{x}_{\text{truth}} + (1 - a)(\underbrace{x_A}_{\text{smoothing error}} - \underbrace{x}_{\text{obs. error}}) + g\varepsilon_o$$

$$a = gk = \frac{\sigma_A^2}{\sigma_A^2 + (\sigma_o/k)^2}$$

averaging kernel

How to solve this inverse problem using Bayesian MAP estimation?

Instead of a single measurement y , make m measurements $\mathbf{y} = (y_1, \dots, y_m)^T$

Cost function:

$$J(x) = \frac{(x - x_A)^2}{\sigma_A^2} + \sum_{i=1}^m \frac{(y_i - kx)^2}{\sigma_o^2} = \frac{(x - x_A)^2}{\sigma_A^2} + \overbrace{\frac{(y_i - kx)^2}{\sigma_o^2 / m}}^{\text{mean}} \quad \text{VS.} \quad J(x) = \frac{(x - x_A)^2}{\sigma_A^2} + \frac{(y - kx)^2}{\sigma_o^2}$$

Observational error variance decreases as m (central limit theorem) IF:

- Error is random
- Individual measurements are not correlated and sample the same pdf

If there is systematic bias in observing system, the bias will propagate to the solution:

$$y = \underbrace{kx}_{\text{bias}} + \underbrace{b_o}_{\text{residual error}} + \varepsilon'_o \quad \Rightarrow \quad \hat{x} = \underbrace{x}_{\text{truth}} + \underbrace{(1-a)(x_A - x)}_{\text{smoothing error}} + \underbrace{g\varepsilon'_o}_{\text{obs. error}} + \underbrace{gb_o}_{\text{bias}}$$

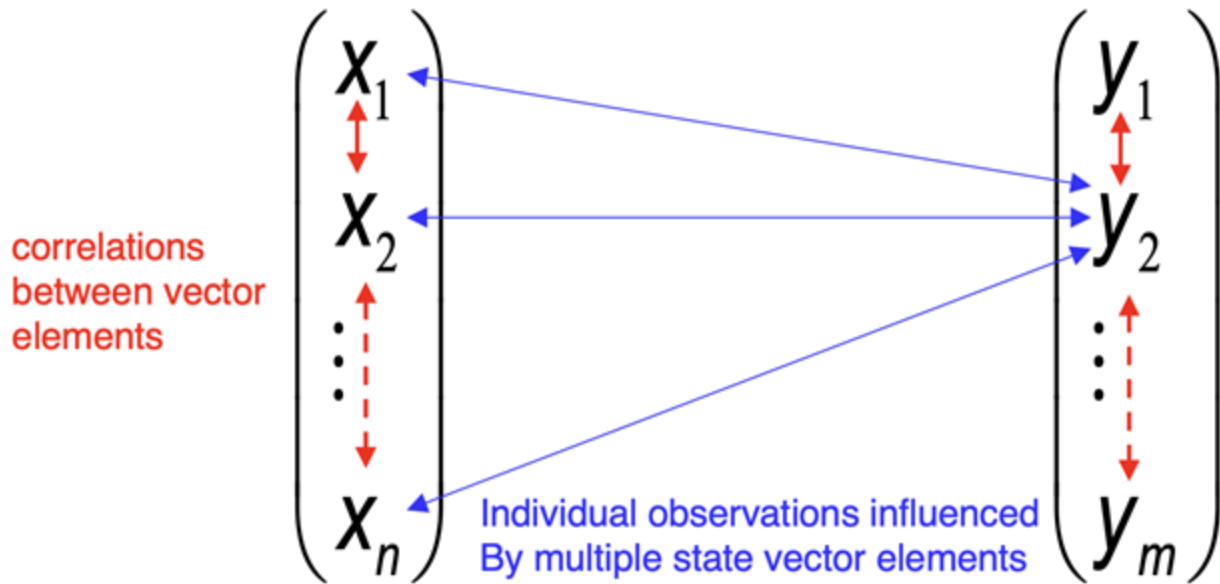
Removing bias by validation of the observing system is critical!

How to formulate this inverse problem using vectors?

Use m observations to constrain a n -dimensional state vector

State vector \mathbf{X} (dimension n)

Observation vector \mathbf{y} (dimension m)



How to formulate error variance matrices and PDFs for the vectors?

Consider vector \mathbf{x} with expected value $E[\mathbf{x}]$ and error $\boldsymbol{\varepsilon} = (\varepsilon_1, \varepsilon_2, \dots, \varepsilon_n)^T$

Error covariance matrix \mathbf{S} :

$$\mathbf{S} = \begin{pmatrix} \text{var}(\varepsilon_1) & \dots & \text{cov}(\varepsilon_1, \varepsilon_n) \\ \vdots & \ddots & \vdots \\ \text{cov}(\varepsilon_1, \varepsilon_n) & \dots & \text{var}(\varepsilon_n) \end{pmatrix}$$

Gaussian error PDF:

$$P(\mathbf{x}) = \frac{1}{(2\pi)^{n/2} |\mathbf{S}|^{1/2}} \exp\left[-\frac{1}{2}(\mathbf{x} - \boxed{})^T \mathbf{S}^{-1}(\mathbf{x} - \boxed{})\right]$$

We will make use of

$$\ln P(\mathbf{x}) \propto -(\mathbf{x} - \boxed{})^T \mathbf{S}^{-1}(\mathbf{x} - \boxed{})$$

How to formulate error variance matrices and PDFs for the vectors?

Consider vector \mathbf{x} with expected value $E[\mathbf{x}]$ and error $\boldsymbol{\varepsilon} = (\varepsilon_1, \varepsilon_2, \dots, \varepsilon_n)^T$

Error covariance matrix \mathbf{S} :

$$\mathbf{S} = \begin{pmatrix} \text{var}(\varepsilon_1) & \dots & \text{cov}(\varepsilon_1, \varepsilon_n) \\ \vdots & \ddots & \vdots \\ \text{cov}(\varepsilon_1, \varepsilon_n) & \dots & \text{var}(\varepsilon_n) \end{pmatrix}$$

Gaussian error PDF:

$$P(\mathbf{x}) = \frac{1}{(2\pi)^{n/2} |\mathbf{S}|^{1/2}} \exp\left[-\frac{1}{2}(\mathbf{x} - E[\mathbf{x}])^T \mathbf{S}^{-1}(\mathbf{x} - E[\mathbf{x}])\right]$$

We will make use of

$$\ln P(\mathbf{x}) \propto -(\mathbf{x} - E[\mathbf{x}])^T \mathbf{S}^{-1}(\mathbf{x} - E[\mathbf{x}])$$

How to formulate Bayesian MAP estimation using vectors?

State vector $\mathbf{x} = (x_1, \dots, x_n)^T$ Obs vector $\mathbf{y} = (y_1, \dots, y_m)^T$ Forward model $\mathbf{y} = \mathbf{F}(\mathbf{x})$

Prior \mathbf{x}_A
with prior error covariance matrix \mathbf{S}_A

Observation \mathbf{y}
with observational error covariance matrix \mathbf{S}_O

prior pdf $P(\mathbf{x})$
 $\ln P(\mathbf{x}) \propto -(\mathbf{x} - \mathbf{x}_A)^T \mathbf{S}_A^{-1} (\mathbf{x} - \mathbf{x}_A)$

observational pdf $P(\mathbf{y} | \mathbf{x})$
 $\ln P(\mathbf{y} | \mathbf{x}) \propto -(\mathbf{y} - \mathbf{F}(\mathbf{x}))^T \mathbf{S}_O^{-1} (\mathbf{y} - \mathbf{F}(\mathbf{x}))$

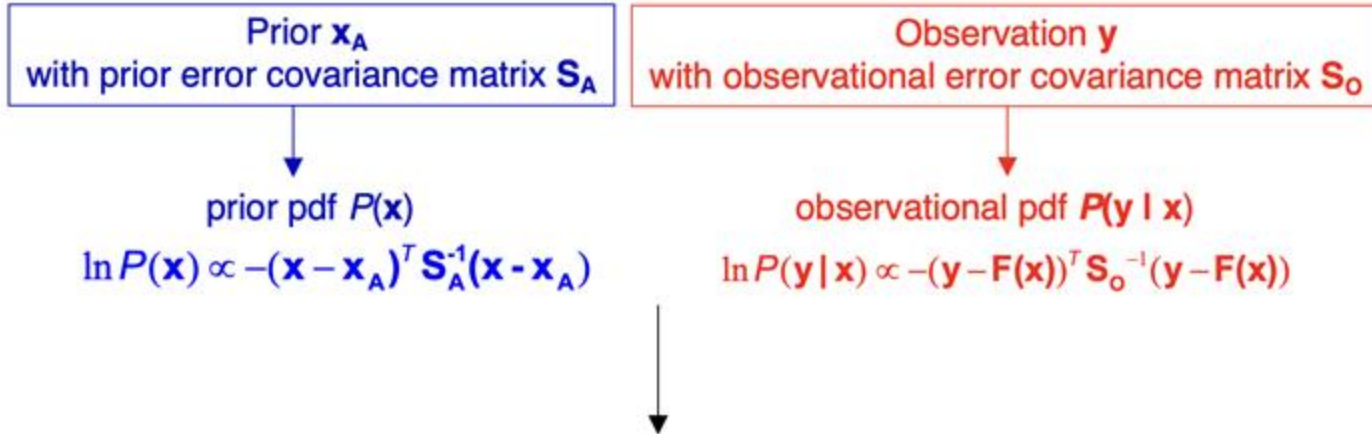
$\ln P(\mathbf{x} | \mathbf{y}) \propto$ $+$ $\ln P(\mathbf{y} | \mathbf{x}) =$ $- (\mathbf{y} - \mathbf{F}(\mathbf{x}))^T \mathbf{S}_O^{-1} (\mathbf{y} - \mathbf{F}(\mathbf{x}))$

Minimize cost function: $J(\mathbf{x}) =$ $+ (\mathbf{y} - \mathbf{F}(\mathbf{x}))^T \mathbf{S}_O^{-1} (\mathbf{y} - \mathbf{F}(\mathbf{x}))$

\Rightarrow solve $\nabla_{\mathbf{x}} J(\mathbf{x}) = 2\mathbf{S}_A^{-1}(\mathbf{x} - \mathbf{x}_A) + 2(\nabla_{\mathbf{x}} \mathbf{F})^T \mathbf{S}_O^{-1} (\mathbf{F}(\mathbf{x}) - \mathbf{y}) = \mathbf{0}$

How to formulate Bayesian MAP estimation using vectors?

State vector $\mathbf{x} = (x_1, \dots, x_n)^T$ Obs vector $\mathbf{y} = (y_1, \dots, y_m)^T$ Forward model $\mathbf{y} = \mathbf{F}(\mathbf{x})$



$$\ln P(\mathbf{x} | \mathbf{y}) \propto \ln P(\mathbf{x}) + \ln P(\mathbf{y} | \mathbf{x}) = -(\mathbf{x} - \mathbf{x}_A)^T \mathbf{S}_A^{-1} (\mathbf{x} - \mathbf{x}_A) - (\mathbf{y} - \mathbf{F}(\mathbf{x}))^T \mathbf{S}_O^{-1} (\mathbf{y} - \mathbf{F}(\mathbf{x}))$$

Minimize cost function: $J(\mathbf{x}) = (\mathbf{x} - \mathbf{x}_A)^T \mathbf{S}_A^{-1} (\mathbf{x} - \mathbf{x}_A) + (\mathbf{y} - \mathbf{F}(\mathbf{x}))^T \mathbf{S}_O^{-1} (\mathbf{y} - \mathbf{F}(\mathbf{x}))$

$$\implies \text{solve } \nabla_{\mathbf{x}} J(\mathbf{x}) = 2\mathbf{S}_A^{-1} (\mathbf{x} - \mathbf{x}_A) + 2(\nabla_{\mathbf{x}} \mathbf{F})^T \mathbf{S}_O^{-1} (\mathbf{F}(\mathbf{x}) - \mathbf{y}) = \mathbf{0}$$

Linear forward model leads to an analytical (closed-form) solution

$$\mathbf{y} = \mathbf{F}(\mathbf{x}) = \mathbf{K}\mathbf{x} \quad \text{where } \mathbf{K} = \begin{pmatrix} \partial y_1 / \partial x_1 & \cdots & \partial y_1 / \partial x_n \\ \vdots & \ddots & \vdots \\ \partial y_m / \partial x_1 & \cdots & \partial y_m / \partial x_n \end{pmatrix} \quad \text{Is the Jacobian matrix}$$

Then equation from previous slide becomes

$$\nabla_{\mathbf{x}} J(\mathbf{x}) = 2\mathbf{S}_A^{-1}(\mathbf{x} - \mathbf{x}_A) + 2\mathbf{K}^T \mathbf{S}_O^{-1}(\mathbf{K}\mathbf{x} - \mathbf{y}) = \mathbf{0}$$

$$\text{Solution: } \hat{\mathbf{x}} = \mathbf{x}_A + \mathbf{G}(\mathbf{y} - \mathbf{K}\mathbf{x}_A) \quad \text{with } \mathbf{G} = \mathbf{S}_A \mathbf{K}^T (\mathbf{K} \mathbf{S}_A \mathbf{K}^T + \mathbf{S}_O)^{-1}$$

$$\hat{\mathbf{S}} = (\mathbf{K}^T \mathbf{S}_O^{-1} \mathbf{K} + \mathbf{S}_A^{-1})^{-1}$$

gain matrix

posterior error covariance matrix

Relate solution to true value:

$$\hat{\mathbf{x}} = \mathbf{x} + (\mathbf{I}_n - \mathbf{A})(\mathbf{x}_A - \mathbf{x}) + \mathbf{G}\boldsymbol{\varepsilon}_O$$

truth

smoothing
error

observational
error

with $\mathbf{A} = \mathbf{G}\mathbf{K}$

averaging kernel matrix

How to calculate K using nonlinear forward model?

Constructing the Jacobian matrix \mathbf{K} of a forward model $\mathbf{y} = \mathbf{F}(\mathbf{x})$

$$\mathbf{K} = \nabla_{\mathbf{x}} \mathbf{F} = \frac{\partial \mathbf{y}}{\partial \mathbf{x}} = \begin{pmatrix} \partial y_1 / \partial x_1 & \dots & \partial y_1 / \partial x_n \\ \vdots & \ddots & \vdots \\ \partial y_m / \partial x_1 & \dots & \partial y_m / \partial x_n \end{pmatrix}$$

Jacobian matrix expresses the sensitivity of elements in \mathbf{y} to elements in \mathbf{x} .

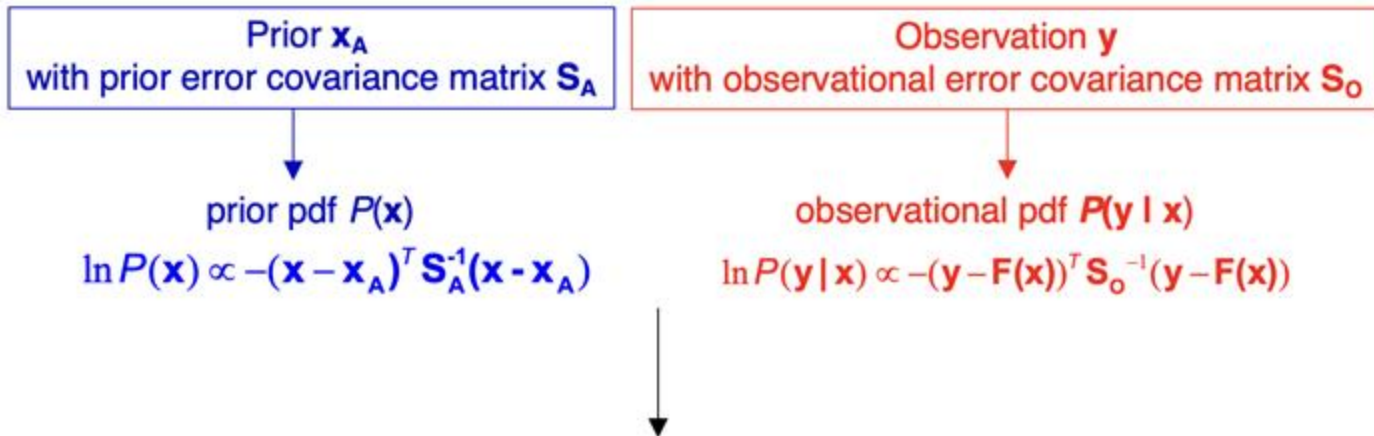
two passes of forward model

$$\frac{\partial \mathbf{y}}{\partial x_1} \approx \frac{\mathbf{F}(\mathbf{x}_A + (\Delta x_1, 0, \dots, 0)^T) - \mathbf{F}(\mathbf{x}_A)}{\Delta x_1}$$

Constructing Jacobian matrix requires $n + 1$ passes of forward model

How to integrate neural network into this framework? Replace \mathbf{x}_A ?

State vector $\mathbf{x} = (x_1, \dots, x_n)^T$ Obs vector $\mathbf{y} = (y_1, \dots, y_m)^T$ Forward model $\mathbf{y} = \mathbf{F}(\mathbf{x})$



$$\ln P(\mathbf{x} | \mathbf{y}) \propto \ln P(\mathbf{x}) + \ln P(\mathbf{y} | \mathbf{x}) = -(\mathbf{x} - \mathbf{x}_A)^T \mathbf{S}_A^{-1} (\mathbf{x} - \mathbf{x}_A) - (\mathbf{y} - \mathbf{F}(\mathbf{x}))^T \mathbf{S}_O^{-1} (\mathbf{y} - \mathbf{F}(\mathbf{x}))$$

Minimize cost function: $J(\mathbf{x}) = (\mathbf{x} - \mathbf{x}_A)^T \mathbf{S}_A^{-1} (\mathbf{x} - \mathbf{x}_A) + (\mathbf{y} - \mathbf{F}(\mathbf{x}))^T \mathbf{S}_O^{-1} (\mathbf{y} - \mathbf{F}(\mathbf{x}))$

$$\implies \text{solve } \nabla_{\mathbf{x}} J(\mathbf{x}) = 2\mathbf{S}_A^{-1} (\mathbf{x} - \mathbf{x}_A) + 2(\nabla_{\mathbf{x}} \mathbf{F})^T \mathbf{S}_O^{-1} (\mathbf{F}(\mathbf{x}) - \mathbf{y}) = \mathbf{0}$$

Questions?