

Linear Regression Introduction

Linear Basis Function Model

1. Least Square Model.

$$y(x, w) = w^T \phi(x)$$

$$w = [w_0, w_1, \dots, w_M]^T \quad \dots \text{parameter} \in \mathbb{R}^M$$

$$\phi(x) = [1, \phi_1(x), \dots, \phi_M(x)]^T \quad \in \mathbb{R}^M$$

w_0 is the bias term.

$\phi_j(x)$ is known as basis function.

For polynomial basis function, $\phi_j(x) = x^j$

For gaussian kernel, $\phi_j(x) = \exp\left(-\frac{(x-\mu_j)^2}{2s^2}\right)$

For sigmoid function, $\phi_j(x) = \sigma\left(\frac{x-\mu_j}{s}\right)$ $\sigma(a) = \frac{1}{1+\exp(-a)}$

Fourier basis

We consider the following additive noise ϵ

$$t = y(x, w) + \epsilon$$

↑ ↑ ↑
received true value $\mathcal{N}(0, \beta^{-1})$
data

β is called precision,
which is the reciprocal of
variance.

So, we can have a series of pdf.

* For input x , the output is t . precision of the noise is β , then.

$$P(t|x, w, \beta) = \mathcal{N}(t | \underbrace{y(x, w)}_{\text{mean}}, \underbrace{\beta^{-1}}_{\text{variance}})$$

if we choose mean-square error. the optimal prediction for t would be $y(x, w)$ with the help of MAP/ML, also.

$$E[t|x] = \int t \cdot P(t|x) dt = y(x, w)$$

* For many input $\mathbf{X} = [x_1, x_2, \dots, x_N]^T$. corresponding output $T = [t_1, \dots, t_N]^T$,

we have joint pdf

$$P(T|\mathbf{X}, w, \beta) = \prod_{i=1}^N \mathcal{N}(t_i | y(x_i, w), \beta^{-1})$$

$$\Rightarrow \ln P(T|\mathbf{X}, w, \beta) = \sum_{i=1}^N \ln \frac{\sqrt{\beta}}{\sqrt{2\pi}} e^{-\frac{1}{2} (t_i - y(x_i, w))^2 \beta}$$

$$= \underbrace{\frac{N}{2} \ln \beta - \frac{N}{2} \ln(2\pi)}_{\text{const}} - \frac{1}{2} \beta \sum_{i=1}^N (t_i - y(x_i, w))^2$$

Now, use ML to compute w, β separately. $\because y(x_i, w) = w^T \phi(x_i)$
 $\because \frac{\partial}{\partial w} y = \phi(x_i)$

** For w .

$$\frac{\partial}{\partial w} \ln P(T|\mathbf{X}, w, \beta) = -\frac{1}{2} \beta \sum_{i=1}^N 2(t_i - y(x_i, w)) \cdot \left(- \frac{\partial}{\partial w} y(x_i, w) \right)$$

$$= \beta \sum_{i=1}^N (t_i - w^T \phi(x_i)) \cdot \phi(x_i)^T$$

Set the gradient to 0

$$0 = \beta \sum_{i=1}^N (t_i - w^T \phi(x_i)) \phi(x_i)^T$$

$$\Rightarrow \sum_{i=1}^N t_i \phi(x_i)^T = w^T \sum_{i=1}^N \phi(x_i) \phi(x_i)^T$$

$$\Leftrightarrow T^T \phi = W^T \phi^T \phi \Leftrightarrow \phi^T T = \phi^T \phi W \Rightarrow W^* = (\phi^T \phi)^{-1} \phi^T T = \phi^\dagger T$$

Note that

$$\phi = \begin{bmatrix} \phi_0(x_1) & \dots & \phi_{M-1}(x_1) \\ \vdots & & \vdots \\ \phi_0(x_N) & & \phi_{M-1}(x_N) \end{bmatrix}_{N \times M}$$

** for β

$$\frac{\partial}{\partial \beta} \ln P(T | \underline{x}, W, \beta) = \frac{N}{2} \cdot \frac{1}{\beta} - \frac{1}{2} \sum_{i=1}^N (t_i - W^*{}^T \phi(x_i))^2 = 0$$

$$\Rightarrow \frac{1}{\beta} = \frac{1}{N} \sum_{i=1}^N (t_i - W^*{}^T \phi(x_i))^2$$

* The error

$$E(W) \equiv \frac{1}{2} \sum_{i=1}^N (t_i - W^T \phi(x_i))^2$$

Definition of
Error

$$= \frac{1}{2} \sum_{i=1}^N \left(t_i - w_0 - \sum_{j=1}^{M-1} w_j \phi_j(x_i) \right)^2$$

$$\frac{\partial}{\partial w_0} E = - \sum_{i=1}^N (t_i - w_0 - \sum_{j=1}^{M-1} w_j \phi_j(x_i)) = 0$$

$$\Rightarrow w_0 = \frac{1}{N} \sum_{i=1}^N t_i - \frac{1}{N} \sum_{i=1}^N \sum_{j=1}^{M-1} w_j \phi_j(x_i)$$

↑
the bias

Geometry Interpretation of least square solution

Our target is to find a parameter W which will minimize the error function:

$$\min_W \frac{1}{2} \sum_{i=1}^N (t_i - W^T \phi(x_i))^2$$

rewritten in matrix form

$$\min_W \frac{1}{2} \|T - \Phi(X)W\|_2^2$$

We can have the following interpretation. (two types.)

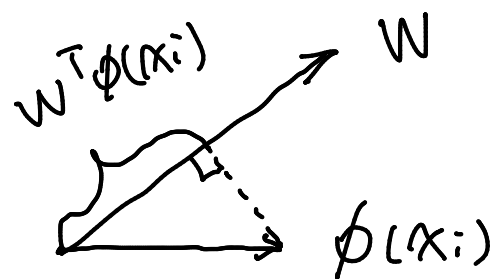
① $T \in \mathbb{R}^N$ target projection distance for N points.

$W \in \mathbb{R}^M$ is a normal vector of $M-D$ plane.

$\Phi(X) \in \mathbb{R}^{M \times N}$ is a set of N points.

So, $W^T \phi(x_i)$ is the projection distance of $M-D$ vector $\phi(x_i)$

Our target is to find a normal, or a direction which will minimize the difference of projection distance and target distance.



② $T \in \mathbb{R}^N$ is a point in $N-D$ space

$W \in \mathbb{R}^M$ is a normal vector of a $M-D$ plane

$\Phi \in \mathbb{R}^{M \times N}$ are N points on $M-D$ plane

Our target is to find the projection w such that the projection

of T onto the M -D space is orthogonal to the N -D plane.

Suppose the projection points

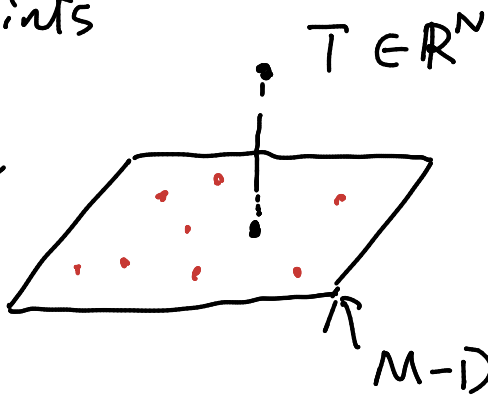
of T on the M -D plane is T' .

Then $\|T - T'\|_2^2$ is the error.

We know that, when TT' is

perpendicular to the M -D plane, $\|T - T'\|_2$ is the smallest. In mathematical

words, we should use "orthogonal" instead of "perpendicular".



• points in M -D space.

Sequential Updating Method * Gradient Descent

Sequential Updating Method is an online algorithm, which is very useful when the data is extremely large.

Reasoning

From previous derivation, we know that the analytical solution is

$$W = \Phi^+ T$$

But the complexity for pseudo inverse is $O(n^3)$, so usually computing Φ^+ brutally is impractical. So, we think of another method. Gradient descent.

Gradient Descent has performance guarantee for convex problem.

But for non-convex problem, no guarantee. but usually have good performance.

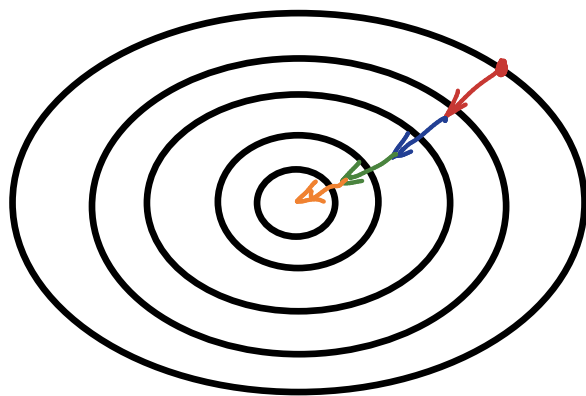
For Convex Problem, local maximum \Leftrightarrow global maximum.

What Gradient Descent does is to go towards local maximum. So, we will always heading to the destination after iterations.

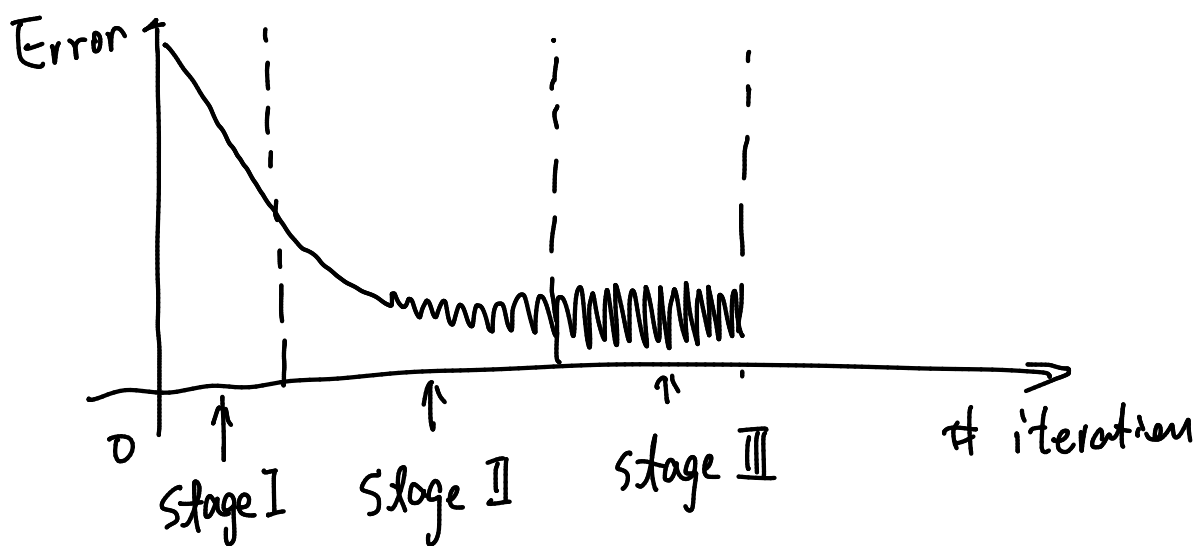
①

$$W = W - \eta \nabla_w E$$

↑
step size.



- $\nabla_w E$ is the gradient updating direction. Each time, we make progress to that direction by step η . Usually, the error would be



Stage I:

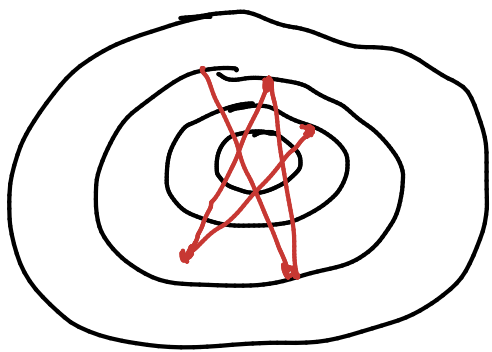
step size η is not large enough. Each update is making roughly linear improvement.

Stage II.

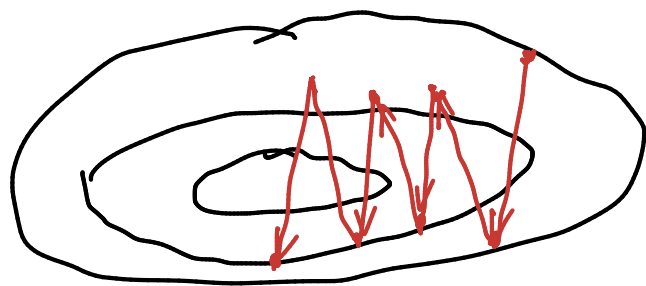
Two explanations

① We are approaching the destination. But η is a little bit larger. We are kind of rambling around the dest.

② We are approaching destination, but the gradient is not quite align with "true" direction.



Case I



Case II

We can choose to decrease step size η to cope with case I, but for case II, Momentum Method may be better. Of course, Newton's method is the optimal solution.

Stage III

Each time, we are working cross the true point. So the performance is fluctuating. We should decrease step size η or terminate the iteration.

2. Regularized Least Square.

Sometimes, we know ahead of time that real coefficient W doesn't have weird parts (say big number), or it is sparse (only parts of W 's coefficient are non-zero). But simply from previous algorithm, we could not have such constrain.

We want the cost function to consider something of W .

We have the following model.

$$\min_W \underbrace{\frac{1}{2} \|T - \Phi W\|_2^2}_{\text{data-dependent error}} + \underbrace{\lambda f(W)}_{\substack{\text{regularization} \\ \text{coefficient}}} \rightarrow \text{regularization term.}$$

We can have many choices for the regularization term.

Usually, we use p -norm $\|W\|_p$, especially 2-norm, 1-norm.

$$\|W\|_p = \sqrt[p]{\sum_{i=1}^N |W_i|^p} \quad p > 0$$

1-norm

$$\|W\|_1 = |W| = \sum_{i=1}^N |W_i|$$

2-norm

$$\|W\|_2 = \|W\| = \sqrt{W^T W}$$

∞ -norm

$$\|W\|_\infty = \max_i (W_i)$$

Same as before, we start from 2-norm's regularization term.

$$\min_W \frac{1}{2} \sum_{i=1}^N [t_i - w^T \phi(x_i)]^2 + \frac{1}{2} \lambda w^T w$$

rewritten the function to matrix multiplication form.

$$\min_W \frac{1}{2} \|T - \Phi W\|_2^2 + \frac{1}{2} \lambda W^T W$$

$$\Leftrightarrow \min_W (T - \Phi W)^T (T - \Phi W) + \lambda W^T W$$

$$\Leftrightarrow \min_W \underbrace{T^T T}_{\text{no } W} - \underbrace{T^T \Phi W}_{\text{Scalar}} - \underbrace{W^T \Phi^T T}_{\text{Scalar}} + \underbrace{W^T \Phi^T \Phi W + \lambda W^T W}_{\substack{\uparrow \\ \lambda W^T I W \\ \downarrow \\ \text{merge}}}$$

$$\Leftrightarrow \min_W \underbrace{-2 T^T \Phi W + W^T (\Phi^T \Phi + \lambda I) W}_{f(W)}$$

$$\frac{\partial}{\partial W} f(W) = -2 (T^T \Phi)^T + 2 (\Phi^T \Phi + \lambda I) W = 0$$

$$\Rightarrow (\Phi^T \Phi + \lambda I) W = \Phi^T T \Rightarrow W = (\Phi^T \Phi + \lambda I)^{-1} \Phi^T T$$

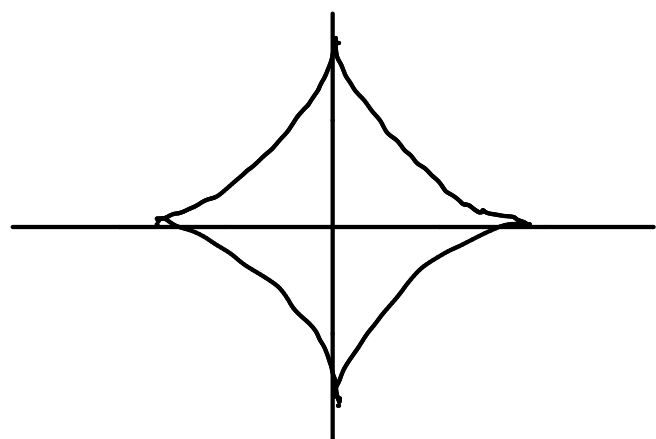
inverse!

PS: $\Phi^T \Phi + \lambda I$ is PSD, so we can do inverse instead of pseudo inverse.

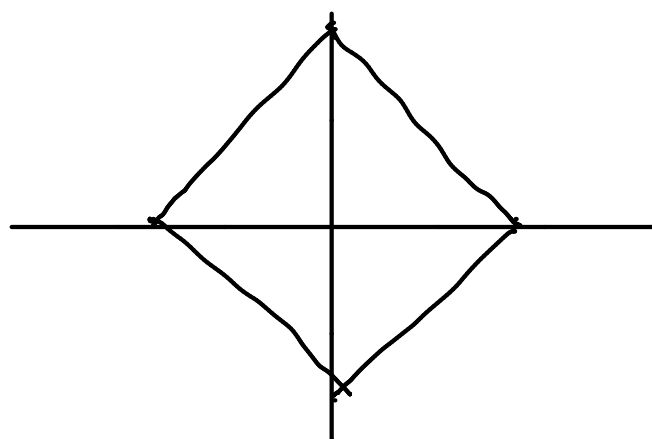
$$\begin{aligned} \text{because } \forall x \in \mathbb{R}^M \quad x^T (\Phi^T \Phi + \lambda I) x &= x^T \Phi^T \Phi x + \lambda x^T x \\ &= \|\Phi x\|_2^2 + \lambda \|x\|_2^2 \geq 0 \end{aligned}$$

What's the difference between different norm?

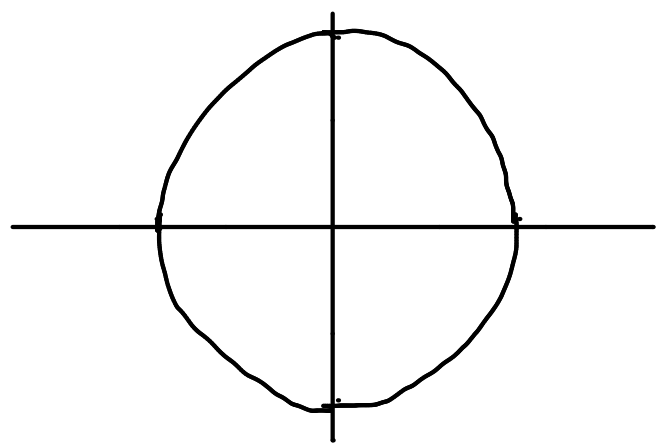
Now, let's see the feasible region for 0.5-norm, 1-norm, 2-norm and 4-norm



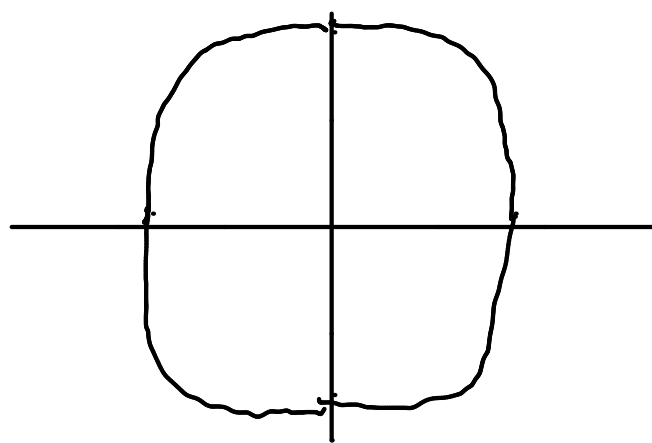
0.5-norm



1-norm



2-norm



4-norm

①

From the figures, we know that 0.5-norm is non-convex.

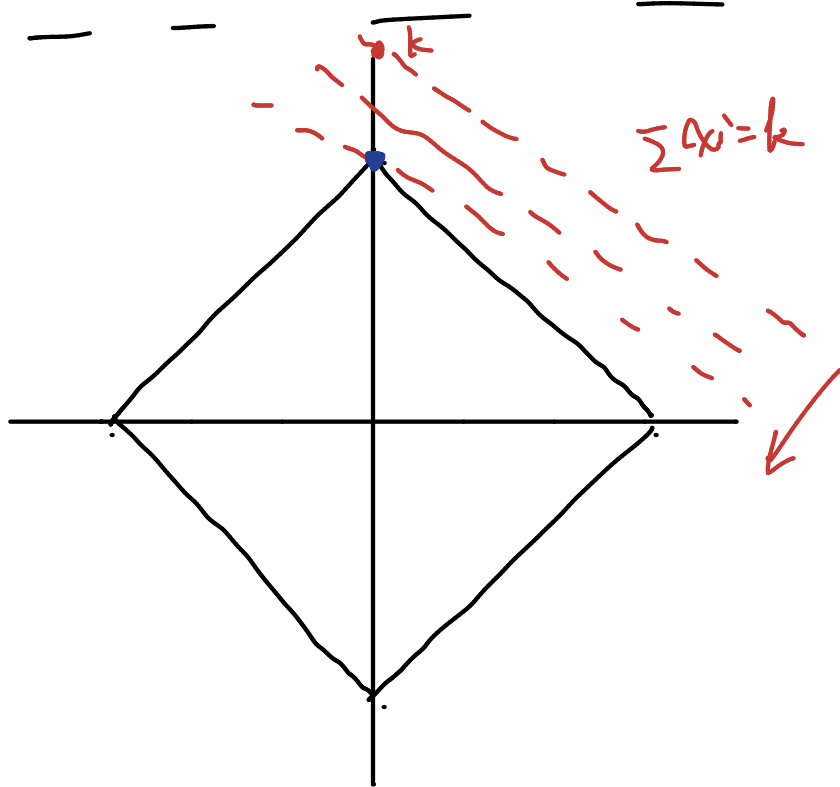
②

If we use 0.5-norm, 1-norm, we are prone to have sparse output.

let's see an linear programming problem

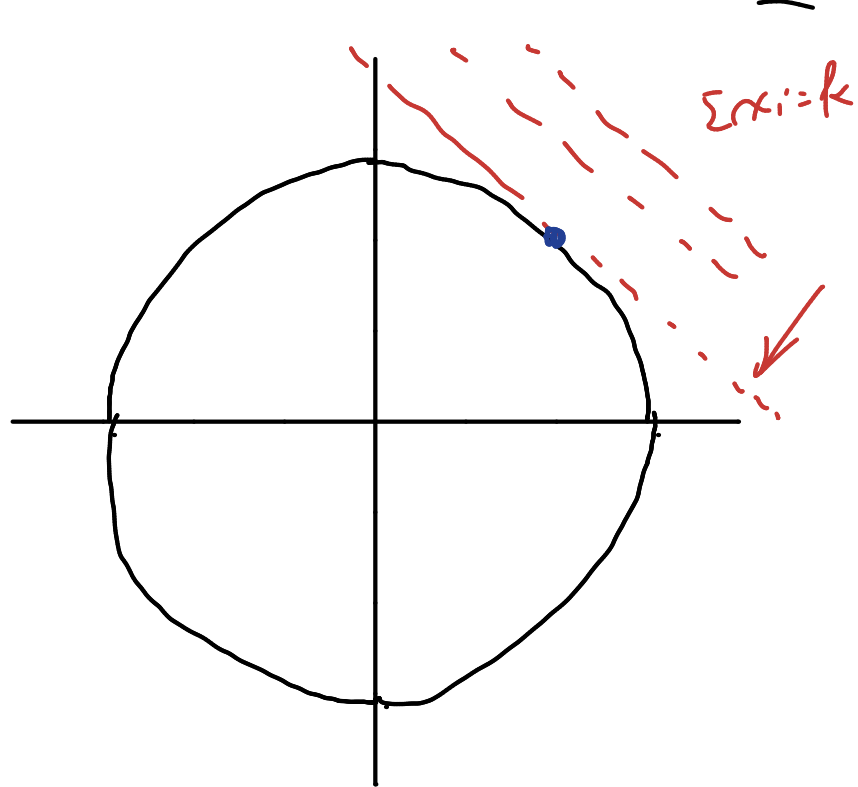
$$\min \sum x_i$$

$$\text{s.t. } |x| \leq 1$$



$$\min \sum x_i$$

$$\text{s.t. } \|x\| \leq 1$$



For 1-norm situation, we will usually have a final point on one axis, so, the final result is sparse, because many other components are 0

For 2-norm, we won't always have such result.

Typically, for 1-norm case, the method is called lasso.

3. The Bias - Variance Decomposition

If we set our model to be extremely complicated, we may easily fall to over-fitting. If we restrict the complexity, we then limit the flexibility of the model. Regularization term is a good way to balance model complexity and overfitting, but how to set coefficient λ ?

Recall that, if our loss function is the "squared" loss function

$$L(x) = x^2$$

the optimal estimator would be the conditional expectation

$$h(x) = \mathbb{E}[t | X]$$

X is the observation, t is the parameter you want to estimate (truth)
 $h(x)$ is the optimal predictor

In real application, we can never know $h(x)$ because we don't know the true pdf behind. So, what we can do is to use some other prediction model $y(x)$. Then, we can have our loss

$$\mathbb{E}[L] = \underbrace{\int (y(x) - h(x))^2 p(x) dx}_{(I)} + \underbrace{\int (h(x) - t)^2 p(x, t) dx dt}_{(II)}$$

part (I) is the difference between our estimator and optimal one.

This term is always non-negative. If we have sufficiently large dataset, we can make this number small.

part (II) is the error we can never overcome, as you can see, there is nothing we can control. This is the noise.

In reality, what we have is a dataset D which is sampled from $p(x, t)$.

For any given dataset D , we can run the prediction function and have $y(x; D)$

Different dataset will give different $y(x; D)$.

Now, part (I) will turn to

$$\int (y(x; D) - h(x))^2 p(x) dx$$

$$\begin{aligned} & (y(x; D) - h(x))^2 \\ &= (y(x; D) - \mathbb{E}_D[y(x; D)] + \mathbb{E}_D[y(x; D)] - h(x))^2 \\ &= (y(x; D) - \mathbb{E}_D[y(x; D)])^2 + [\mathbb{E}_D[y(x; D)] - h(x)]^2 \\ & \quad + 2 \underbrace{(y(x; D) - \mathbb{E}_D[y(x; D)])(\mathbb{E}_D[y(x; D)] - h(x))}_{=0} \end{aligned}$$

Now, we take expectation, note that

$$\begin{aligned} & \mathbb{E}_D [2 (y(x; D) - \mathbb{E}_D[y(x; D)])(\mathbb{E}_D[y(x; D)] - h(x))] \\ &= 2(\mathbb{E}_D[y(x; D)] - h(x)) \cdot \underbrace{\mathbb{E}_D [y(x; D) - \mathbb{E}_D[y(x; D)]]}_{=0} \\ &= 0 \end{aligned}$$

$$\begin{aligned} \therefore \mathbb{E}_D [(y(x; D) - h(x))^2] \\ &= \underbrace{\mathbb{E}_D [(y(x; D) - \mathbb{E}_D[y(x; D)])^2]}_{\text{Variance}} + \underbrace{[\mathbb{E}_D[y(x; D)] - h(x)]^2}_{(\text{bias})^2} \end{aligned}$$

So, we can split the expected loss to

expected loss = (bias)² + variance + noise,

$$(\text{bias})^2 = \int (\mathbb{E}_D [y(x; D)] - h(x))^2 p(x) dx$$

$$\text{variance} = \int \mathbb{E}_D [(y(x; D) - \mathbb{E}_D [y(x; D)])^2] p(x) dx$$

$$\text{noise} = \int (h(x) - t)^2 p(x, t) dx dt$$

Remind that, our target is to minimize the overall loss, i.e. minimize the expected loss.

Noise is unsolvable. it's there, and we cannot deal with that. But there is a trade-off between bias - variance.