

# 04-Linear Regression (Probabilistic View)

## Probabilistic assumption

Assume the true value  $t$  (groundtruth) and the input  $x$  satisfies the following equation

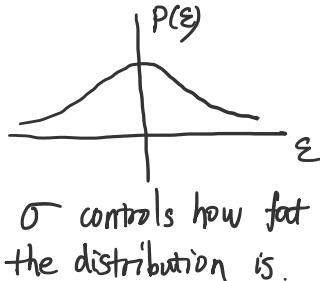
$$t = \mathbf{w}^T \mathbf{x} + \varepsilon$$

where  $\varepsilon \sim N(0, \sigma^2)$

That is to say,  $t$  is indeed a linear/affine transformation of  $x$ , except that some Gaussian noise affects  $t$  in an additive way.

The normal distribution  $N(0, \sigma^2)$ :

$$p(\varepsilon) = \frac{1}{\sqrt{2\pi}\sigma} \exp\left\{-\frac{1}{2}\left(\frac{\varepsilon}{\sigma}\right)^2\right\}$$



It is required to know Gaussian is quadratic in exponential decay.  
Memorizing the normalizing factor is not required in exams.

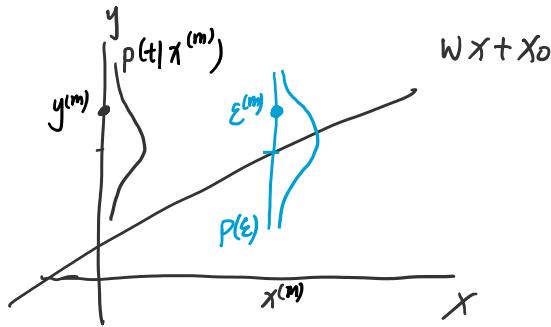
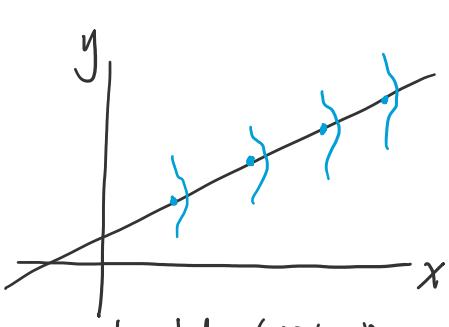
Since we know  $\varepsilon = t - \mathbf{w}^T \mathbf{x}$

$$p(\varepsilon) = \frac{1}{\sqrt{2\pi}\sigma} \exp\left\{-\frac{1}{2}\left(\frac{t - \mathbf{w}^T \mathbf{x}}{\sigma}\right)^2\right\}$$

↑  
The probability density of  $\varepsilon$  taking some value.

↑  
The probability density of  $t$  taking some corresponding value given  $\mathbf{x}$ ,

$$p(t|\mathbf{x}; \mathbf{w}) = \frac{1}{\sqrt{2\pi}\sigma} \exp\left\{-\frac{1}{2}\left(\frac{t - \mathbf{w}^T \mathbf{x}}{\sigma}\right)^2\right\}$$



Knowledge/assumption

A data sample

## Maximum likelihood estimation

Suppose we have a set of data  $\mathcal{D}$

We also have some probabilistic modeling of data

$$P(\mathcal{D}; \theta)$$

where  $\theta$  is the parameter.

Likelihood of parameters is yet another terminology of data probability

$$\mathcal{L}(\theta; \mathcal{D})$$

$$\stackrel{\Delta}{=}$$

$$P(\mathcal{D}; \theta)$$

Maximum likelihood estimation (MLE) maximizes the likelihood of parameters, which is the probability of data.

$$\hat{\theta}_{\text{MLE}} = \arg \max_{\theta} \mathcal{L}(\theta)$$

$$= \arg \max_{\theta} \log \mathcal{L}(\theta)$$

[ $\log$  is monotonically increasing]

$\log$  of likelihood is known as

$\log$ -likelihood, denoted by  $l(\theta; \mathcal{D})$

In our probabilistic assumption for linear regression, we consider the probability of all labels  $t^{(i)}$  given  $x^{(i)}$ , parameterized by model parameters.

$$\begin{aligned} \log \mathcal{L}(w) &= \log P(t^{(1)}, t^{(2)}, \dots, t^{(M)} | x^{(1)}, x^{(2)}, \dots, x^{(M)}) \quad [\text{y's are RVs, x's not}] \\ &= \log \prod_{m=1}^M p(t^{(m)} | x^{(1)}, x^{(2)}, \dots, x^{(M)}) \quad [\text{data iid}] \\ &= \log \prod_{m=1}^M p(t^{(m)} | x^{(m)}) \quad [\text{y}^{(i)} \text{ only depends on } x^{(i)}] \\ &= \sum_{m=1}^M \log p(t^{(m)} | x^{(m)}) \quad [\log \prod = \sum \log] \\ &= \sum_{m=1}^M \log \frac{1}{\sqrt{2\pi}\sigma} \exp \left\{ -\frac{1}{2} \left( \frac{t^{(m)} - w^T x^{(m)}}{\sigma} \right)^2 \right\} \\ &= \sum_{m=1}^M \left[ \log \frac{1}{\sqrt{2\pi}\sigma} + \log \exp \left\{ -\frac{1}{2} \left( \frac{t^{(m)} - w^T x^{(m)}}{\sigma} \right)^2 \right\} \right] \\ &= -M \cdot \log(\sqrt{2\pi}\sigma) - \frac{1}{2} \sum_{m=1}^M \left( \frac{t^{(m)} - w^T x^{(m)}}{\sigma} \right)^2 \end{aligned}$$

Thus,

maximize  $\log \mathcal{L}(w)$

$$\Leftrightarrow \underset{\mathbf{w}}{\text{minimize}} \sum_{m=1}^M (t^{(m)} - \mathbf{w}^T \mathbf{x}^{(m)})^2$$

$$\Leftrightarrow \underset{\mathbf{w}}{\text{minimize}} \frac{1}{2M} \sum_{m=1}^M (t^{(m)} - \mathbf{w}^T \mathbf{x}^{(m)})^2$$

This is exactly MSE.

## Best Linear Unbiased Estimate (BLUE)

**Theorem.** Assume  $t^{(i)}$  is generated by  $t^{(i)} = \mathbf{w}^T \mathbf{x} + \epsilon^{(i)}$ , where  $\epsilon^{(i)}$  is zero-mean, uncorrelated, and with finite variance;  $\mathbf{w}$  is unknown parameters. Then, least square error yields an unbiased estimate of  $\mathbf{w}$ .

Proof: The closed-form solution is  $\hat{\mathbf{w}} = (\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T \mathbf{t}$

$$\begin{aligned} \mathbb{E}_{\varepsilon^{(m)}}[\hat{\mathbf{w}}] &= \mathbb{E}_{\varepsilon^{(m)}}[(\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T \mathbf{t}] \\ &= \mathbb{E}_{\varepsilon^{(m)}}[(\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T (\mathbf{X}\mathbf{w} + \varepsilon)] \quad [\varepsilon = (\varepsilon^{(1)}, \dots, \varepsilon^{(M)})^T] \\ &= \mathbb{E}_{\varepsilon^{(m)}}[(\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T \mathbf{X}\mathbf{w}] + \mathbb{E}_{\varepsilon^{(m)}}[(\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T \varepsilon] \\ &= \mathbf{w} \quad \# \end{aligned}$$

**Theorem (Gauss-Markov).** Least error regression yields the least variance in linear, unbiased estimators. (Here, a linear estimator means  $\mathbf{w} = \mathbf{C}\mathbf{t}$  for some  $\mathbf{C}$ .)

[#](https://en.wikipedia.org/wiki/Gauss-Markov_theorem)

## Bias-Variance Tradeoff

- Limitation of the above analysis:

In reality the assumption  $\mathbf{t} = \mathbf{w}^T \mathbf{x} + \varepsilon$  is usually not true, and we do not know the true function.

Now suppose  $\mathbf{t} = f(\mathbf{x}) + \varepsilon$  where  $\varepsilon \sim p(\varepsilon|\mathbf{x})$  and  $\mathbb{E}_{\varepsilon \sim p(\varepsilon|\mathbf{x})}[\varepsilon] = 0 \quad \forall \mathbf{x}$

In other words, we can treat  $\mathbf{t}$  as an RV. and define  $f(\mathbf{x}) = \mathbb{E}_{\substack{\mathbf{t} \sim p(\mathbf{t}|\mathbf{x})}}[\mathbf{t}]$

We consider the performance of a machine learning model  $h(\mathbf{x})$  by the measure of success. In our case, the measure of success is modeled by MSE. We call how bad the model performs on a data sample  $\mathbf{x}$  the **error** on  $\mathbf{x}$ , denoted by  $E(\mathbf{x})$ . In here, it happens to be the same as the training objective function. We call the expected loss of a model  $h$  (when data are repeatedly generated from the data distribution) **error**, too, denoted by  $E(h)$ . This is also known as the **risk**.

$$\begin{aligned}
 E(h) &= \mathbb{E}_{\mathbf{x}, \varepsilon} [(f(\mathbf{x}) + \varepsilon - h(\mathbf{x}))^2] && \left[ \begin{array}{l} \text{Recall: } h(\mathbf{x}) \text{ is the prediction} \\ f(\mathbf{x}) \text{ is the "true" function,} \\ \text{less noise} \end{array} \right] \\
 &= \mathbb{E}_{\mathbf{x}, \varepsilon} [(f(\mathbf{x}) - h(\mathbf{x}) + \varepsilon)^2] \\
 &= \mathbb{E}_{\mathbf{x}, \varepsilon} [(f(\mathbf{x}) - h(\mathbf{x}))^2] + \mathbb{E}_{\mathbf{x}, \varepsilon} [(f(\mathbf{x}) - h(\mathbf{x})) \cdot \varepsilon] + \mathbb{E}_{\mathbf{x}, \varepsilon} [\varepsilon^2] && (\text{E is linear}) \\
 &&= 0, \text{ because } \mathbb{E}[\varepsilon] = 0 \\
 &= \mathbb{E}_{\mathbf{x}} [(f(\mathbf{x}) - h(\mathbf{x}))^2] + \mathbb{E}_{\substack{\mathbf{x} \sim p(\mathbf{x}) \\ \varepsilon \sim p(\varepsilon | \mathbf{x})}} [\varepsilon^2] && (1)
 \end{aligned}$$

Notice that  $h(\mathbf{x})$ , in fact, depends on the training data  $\mathcal{D}$ . We denote it by  $h_{\mathcal{D}}(\mathbf{x})$ . The first term can be further decomposed as

$$\begin{aligned}
 \mathbb{E}_{\mathbf{x}} [(f(\mathbf{x}) - h_{\mathcal{D}}(\mathbf{x}))^2] &= \mathbb{E}_{\mathbf{x}} [\{f(\mathbf{x}) - \mathbb{E}_{\mathcal{D}}[h_{\mathcal{D}}(\mathbf{x})] + \mathbb{E}_{\mathcal{D}}[h_{\mathcal{D}}(\mathbf{x})] - h_{\mathcal{D}}(\mathbf{x})\}^2] \\
 &\quad \triangleright \text{Did nothing, only introduced } \pm \mathbb{E}_{\mathcal{D}}[h_{\mathcal{D}}(\mathbf{x})] \\
 &= \mathbb{E}_{\mathbf{x}} [\{f(\mathbf{x}) - \mathbb{E}_{\mathcal{D}}[h_{\mathcal{D}}(\mathbf{x})]\}^2] + \mathbb{E}_{\mathbf{x}} [\{\mathbb{E}_{\mathcal{D}}[h_{\mathcal{D}}(\mathbf{x})] - h_{\mathcal{D}}(\mathbf{x})\}^2] \\
 &\quad + \mathbb{E}_{\mathbf{x}} [2(f(\mathbf{x}) - \mathbb{E}_{\mathcal{D}}[h_{\mathcal{D}}(\mathbf{x})]) \cdot (\mathbb{E}_{\mathcal{D}}[h_{\mathcal{D}}(\mathbf{x})] - h_{\mathcal{D}}(\mathbf{x}))] && (2)
 \end{aligned}$$

We further take expectation of data  $\mathcal{D}$  (assuming data can be iid generated from a repeated trial).

In this case  $\mathbb{E}_{\mathcal{D}}[\text{last term above}] = 0$  because  $\mathbb{E}_{\mathcal{D}}[\mathbb{E}_{\mathcal{D}}[h_{\mathcal{D}}(\mathbf{x})] - h_{\mathcal{D}}(\mathbf{x})] = 0$

Combining (1) and (2), we have

$$\mathbb{E}_{\mathcal{D}} [R(h_{\mathcal{D}})] = \underbrace{\mathbb{E}_{\mathbf{x}} [\{f(\mathbf{x}) - \mathbb{E}_{\mathcal{D}}[h_{\mathcal{D}}(\mathbf{x})]\}^2]}_{\text{bias}^2} + \underbrace{\mathbb{E}_{\mathcal{D}, \mathbf{x}} [\{\mathbb{E}_{\mathcal{D}}[h_{\mathcal{D}}(\mathbf{x})] - h_{\mathcal{D}}(\mathbf{x})\}^2]}_{\text{Variance}} + \underbrace{\mathbb{E}_{\substack{\mathbf{x} \sim p(\mathbf{x}) \\ \varepsilon \sim p(\varepsilon | \mathbf{x})}} [\varepsilon^2]}_{\text{noise}}$$

**Interpretation:** The expected error (said in terms of training data  $\mathcal{D}$  repeated drawn from a repeatable trial and test data  $x, \epsilon$  following test distribution) is decomposed of

$$\text{Expected error} = \text{Bias}^2 + \text{Variance} + \text{Noise}$$

- **Bias:** The difference between the true function and the expected learned function
- **Variance:** The square difference between the learned function  $h$  from this training data  $\mathcal{D}$  and the expected learned function
- **Noise:** Variance of noise, which is intrinsic to the dataset and cannot be optimized

#### Reducing bias:

- Increasing the hypothesis class (containing the true function is a necessary condition for being unbiased)

#### Reducing variance:

- Reducing the hypothesis class (fewer functions to choose from)

**Example:** Suppose the truth function is  $n$ th-degree polynomial

$$t = w_1 x + w_2 x^2 + \dots + w_n x^n$$

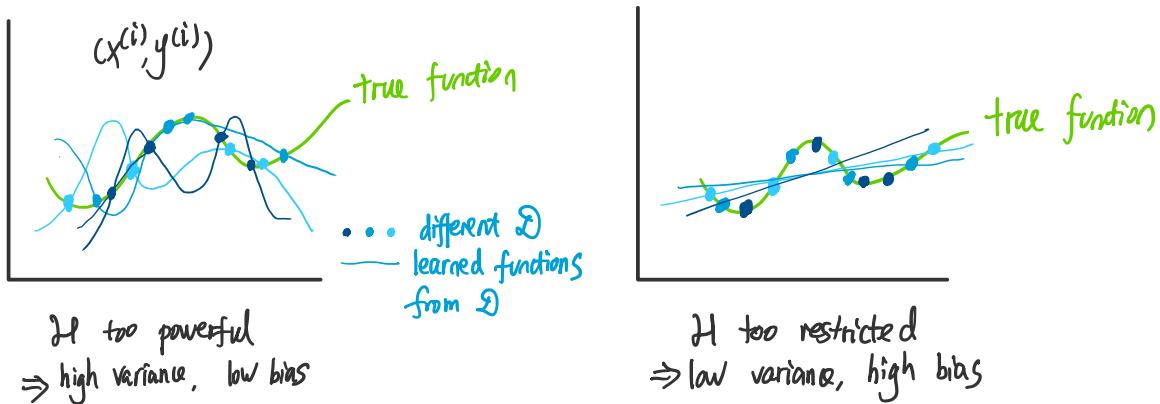
We can introduce polynomial basis as the feature set.

$$x_1 = x, \quad x_2 = x^2, \quad \dots, \quad x_d = x^d, \quad \text{and } x_0 = 1 \text{ as usually.}$$

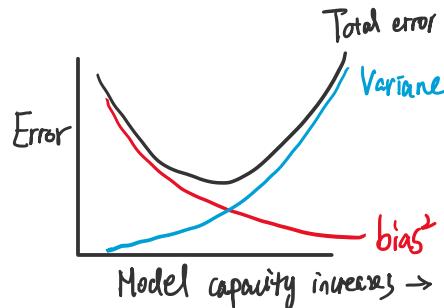
In this case, the least square estimate is still unbiased if  $d \geq n$  (as it has the same format as linear regression). The resulting model is still linear in  $w$  and  $(x_1, \dots, x_d)^T$ , although not linear in  $x$ .

Unfortunately, we do not know the true order of polynomial  $n$ .

- If  $d$  is too large, the model is unbiased, but since we have too many functions to choose from, the estimator is of high variance.
- If, on the other hand,  $d$  is too small, the model is biased, but we have fewer functions to choose from and the estimator is of low variance.



### An intuitive picture of bias-variance tradeoff



### An approach to the bias-variance tradeoff

- We must state preference of  $\mathcal{H}$  before training
- Instead of a highly restricted form of functions (e.g., linear with limited features), we
  - Allow a more powerful form of functions, but
  - State some other preferences, e.g.,
    - Among  $w_1, w_2, \dots, w_d$ , only a few parameters could be non-zero, or
    - The weights  $w_1, w_2, \dots, w_d$  cannot be too large, etc.

### Footnote on the terminologies

- In machine learning, people sometimes use *objective*, *loss*, *risk*, *error* interchangeably. They are indeed similar, but may have different emphasis.
- **Objective**, or **training objective**, is clear. It is the goal of training, which may or may not be the same function as loss/risk/error.
- **Loss** has two meanings.
  - When we say **training loss**, it usually means training objective.
  - ~ In decision theory (where they don't care about training) **loss**

- In decision theory (where they don't care about training), loss means the measure of success for a model  $h$  evaluated on a particular data point  $L(x, h)$ . We use **Error** in this course. The lower, the better.
- **Risk** is usually said by referring to the expected loss/error when data are repeated generated, so it's also called the **frequentist risk**. It is function of your predictor  $h$  and the true model  $f$ , denoted by  $R(h, f)$ . It should be emphasized that the expectation is said in terms of data being repeated generated from the data distribution. So this terminology actually makes some sense, because we have uncertainty of the world (data), and thus, risk.

Reference: Chap 1.3, James O. Berger, *Statistical Decision Theory and Bayesian Analysis*.

<https://link.springer.com/book/10.1007/978-1-4757-4286-2>