## UNITS 1-2

$$y \approx \beta_0 + \beta_1 x_1 + \ldots + \beta_d x_d \overset{\Delta}{=} \hat{y}$$

training :
$$\underbrace{\begin{bmatrix} y_1 \\ \vdots \\ y_n \end{bmatrix}}_{\underline{y}} \approx \underbrace{\begin{bmatrix} 1 & x_{11} & \cdots & x_{1d} \\ \vdots & & & \\ 1 & x_{n1} & \cdots & x_{nd} \end{bmatrix}}_{\underline{A}} \underbrace{\begin{bmatrix} \beta_0 \\ \beta_1 \\ \vdots \\ \beta_d \end{bmatrix}}_{\underline{\beta}} \overset{\Delta}{=} \underbrace{\begin{bmatrix} \hat{y}_1 \\ \vdots \\ \hat{y}_n \end{bmatrix}}_{\underline{\hat{y}}}$$

<span style="color:red">↰ test</span>

test : $\hat{y}(\underline{x}) = \begin{bmatrix} 1 & \underline{x}^T \end{bmatrix} \hat{\underline{\beta}}$ <span style="color:red">← learned</span>

$$RSS(\underline{\beta}) = \sum_{i=1}^{n} (y_i - \hat{y}_i)^2 = \| \underline{y} - \underline{A}\underline{\beta} \|^2 \qquad \text{<span style='color:red'>"least squares"</span>}$$

to find $\hat{\underline{\beta}}_{LS}$ :  $\left. \dfrac{\partial}{\partial \beta_j} RSS(\underline{\beta}) \right|_{\hat{\beta}_j} = 0 \; \forall j \;\Rightarrow\; \hat{\underline{\beta}}_{LS} = (\underline{A}^T \underline{A})^{-1} \underline{A}^T \underline{y}$

$R^2$ coefficient of determination : $R^2 \overset{\Delta}{=} 1 - \dfrac{RSS}{n \, s_y^2}$

<span style="color:red">$\begin{cases} R^2 = 1 : & \text{perfect} \\ R^2 = 0 : & \text{trivial} \\ R^2 < 0 : & \text{impossible on training data for } \hat{\underline{\beta}}_{LS} \end{cases}$</span>

categorical features  $x_j \in \{A, B, C\}$
 - use one-hot coding : turn $x_j$ into a one-hot binary vector
 - coding the intercept : each category gets a unique intercept <span style="color:red">$\beta_0$</span>
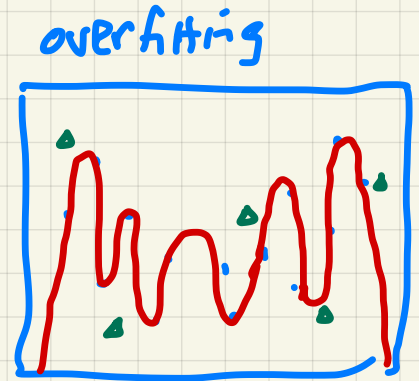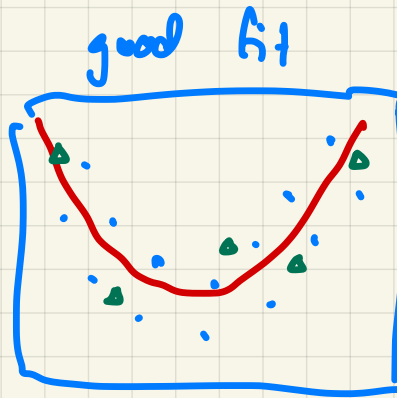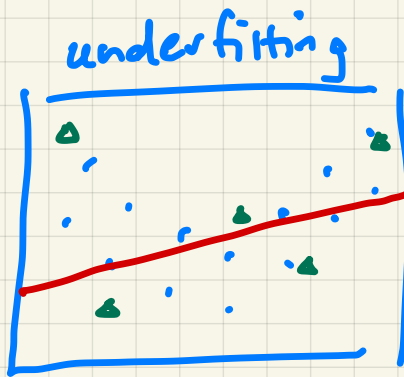 - code the slope :  "  "  "  "  slope

nonlinear transformations ⟶ new features
 e.g. polynomial regression
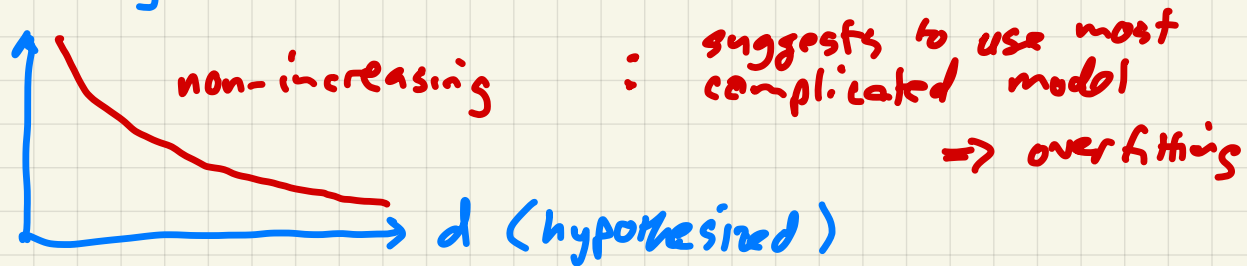$$y \approx \beta_0 + \beta_1 x + \beta_2 x^2 + \ldots + \beta_d x^d \qquad \text{<span style='color:red'>← use } x_j = x^j\text{</span>}$$
 but how to choose d ?

# UNIT 3

underfitting | good fit | overfitting

training RSS (or MSE)

non-increasing : suggests to use most complicated model ⇒ overfitting

$d$ (hypothesized)

solution: cross-validation = choose model order using samples different from training samples (ones used to get $\hat{\beta}$)

options: 1) test / train split

2) K-fold : train on $K-1$ folds, test on remaining fold } repeat for all $K$ combinations

$$\overline{RSS}_d = \frac{1}{K} \sum_{k=1}^{K} RSS_{k,d} \quad \leftarrow \text{estimate of } E\{RSS_d\}$$
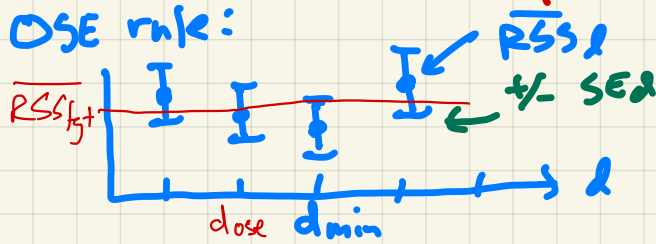
simple approach: choose $d_{min} = \arg\min_d \overline{RSS}_d$ ↑ random

"standard error" on $\overline{RSS}_d$

$$SE_d \triangleq \frac{\hat{\sigma}}{\sqrt{K}} \quad \text{where} \quad \hat{\sigma} = \sqrt{\frac{1}{K-1} \sum_{i=1}^{K} \left( RSS_{i,d} - \overline{RSS}_d \right)^2}$$

↳ estimate of $\sqrt{\text{var}\{RSS_d\}}$

OSE rule:

$\overline{RSS}_d$

$\overline{RSS}_{tgt}$ +/- $SE_d$

dose $d_{min}$

$$\overline{RSS}_{tgt} = \overline{RSS}_{d_{min}} + SE_{d_{min}}$$

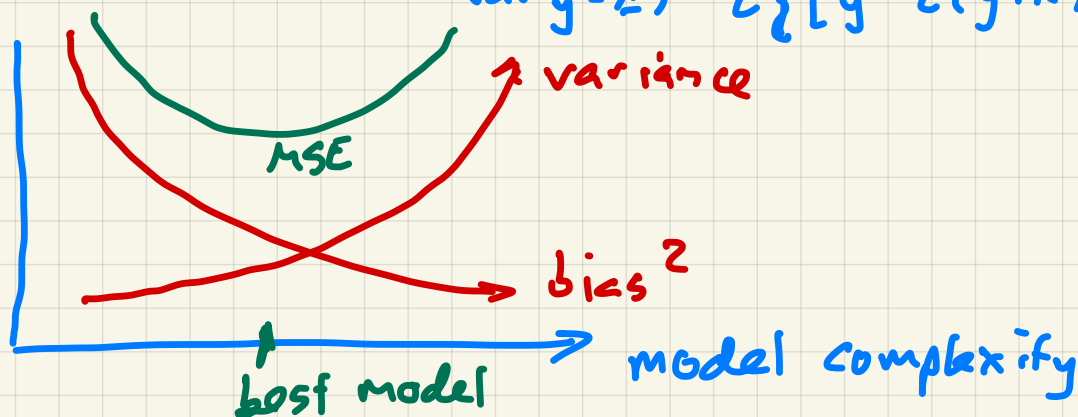$$d_{ose} = \min_d \left\{ d : \overline{RSS}_d \leq \overline{RSS}_{tgt} \right\}$$

$d$ ← simplest

# Bias - Variance tradeoff

- true model: training: $y_i = f(x_i) + \varepsilon_i$ } $\varepsilon_i$ and $\varepsilon$ are
  testing: $y = f(x) + \varepsilon$ } iid, zero-mean, variance $\sigma^2$

- prediction model: $\hat{y} = \hat{f}(x, \hat{\beta})$ ← random trained parameters depend on random $x_i$ & $\varepsilon_i$

- metric: $MSE_{\hat{y}}(x) \triangleq E\{(y-\hat{y})^2 | x\}$

  we derived $MSE_{\hat{y}}(x) = \sigma^2 + \left[bias_{\hat{y}}(x)\right]^2 + var_{\hat{y}}(x)$

  where $bias_{\hat{y}}(x) = E\{\hat{y}-y | x\}$

  $var_{\hat{y}}(x) = E\left\{ \left[\hat{y} - E\{\hat{y}|x\}\right]^2 \Big| x\right\}$



- Special case: LS linear regression

  $d < d_{true}$ : $bias \neq 0$ ⟹ underfitting

  $d \geq d_{true}$ : $bias = 0$

  $E\{var_{\hat{y}}(x)\} = \dfrac{d+1}{n}\sigma^2$

$\boxed{\text{UNIT 4}}$ Feature Selection: choose best subset of $d$ features

- exhaustive search: optimal, but complexity grows as $2^d$
- stepwise regression: greedy, but useful ... complexity $d^2$
- ranking based on univariate statistics (correlation)
- regularization - based methods

<span style="color:red">L1 regularization</span>

① LASSO: $\arg\min_{\underline{\beta}} \{ \|\underline{y} - \underline{X}\underline{\beta}\|_2^2 + \alpha \|\underline{\beta}\|_1 \}$

    · sets a subset of $\{\beta_j\}$ to zero, shrinks remaining $\beta_j$
    · we <u>use</u> the indices of nonzero $\{\beta_j\}$ but not their values

<span style="color:red">fit a LS model</span>

<span style="color:red">L2 regularization</span>

② Ridge: $\arg\min_{\underline{\beta}} \{ \|\underline{y} - \underline{X}\underline{\beta}\|_2^2 + \alpha \|\underline{\beta}\|_2^2 \}$

    · not useful for feature selection ( no zero-valued $\beta_j$ )
    · useful with correlated features

## Probabilistic Interpretations

<span style="color:red">likelihood from fixed</span>

① Maximum Likelihood (ML): $\hat{\underline{\beta}}_{ML} \triangleq \arg\max_{\underline{\beta}} p(\underline{y} | \underline{X}, \underline{\beta})$
$$= \arg\min_{\underline{\beta}} [ -\ln p(\underline{y} | \underline{X}, \underline{\beta}) ]$$

   · ML estimation of $\underline{\beta}$
      under $\underline{y} = \underline{X}\underline{\beta} + \underline{\varepsilon}, \; \underline{\varepsilon} \sim N(\underline{0}, \sigma^2 I)$
      gives LS estimation: $\hat{\underline{\beta}}_{ML} = \arg\min_{\underline{\beta}} \|\underline{y} - \underline{X}\underline{\beta}\|_2^2$

② MAP: $\hat{\underline{\beta}}_{MAP} = \arg\max_{\underline{\beta}} p(\underline{\beta} | \underline{X}, \underline{y})$   <span style="color:red">$p(\underline{\beta}|\underline{X},\underline{y}) = \dfrac{p(\underline{y}|\underline{X},\underline{\beta}) p(\underline{\beta})}{p(\underline{y}|\underline{X})}$</span>
$$= \arg\min_{\underline{\beta}} [ -\ln p(\underline{\beta} | \underline{X}, \underline{y}) ]$$

   · Ridge is MAP under $\beta_j \sim N(0, v)$     } and linear/ Gaussian likelihood
   · LASSO is MAP under $\beta_j \sim$ Laplacian $(v)$