

Linear prediction: features, overfitting, and losses

Feature Engineering:

$\mathbb{R}^d \rightarrow \mathbb{R}^p$ replace with x with $\phi(x)$

Eg ① non-linear transformation $x \in \mathbb{R}$

semi-often $\phi(x) = \ln(1+x)$

② logical formula, for $x = (x_1, \dots, x_d) \in \{0, 1\}^d$

never seen it $\phi(x) = (x_1 \wedge x_5 \wedge \neg x_{10}) \vee (\neg x_2 \wedge x_7)$

③ Trigonometric expansion $x \in \mathbb{R}$

rare $\phi(x) = (1, \sin(x), \cos(x), \sin(2x), \cos(2x), \dots)$

④ polynomial expansion: for $x = (x_1, \dots, x_d) \in \mathbb{R}^d$

frequent $\phi(x) = (1, x_1, \dots, x_d, x_1^2, \dots, x_d^2, x_1 x_2, \dots, x_1 x_d, \dots, x_{d-1} x_d)$

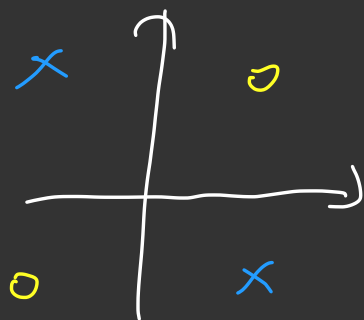
Monomial features

$$\phi(x) = (1, x, x^2) \quad w = (a, b, c)$$

$$x \mapsto w^T \phi(x) = a + bx + cx^2$$

$$w^T \phi(x) = \begin{bmatrix} a \\ b \\ c \end{bmatrix}^T \begin{bmatrix} 1 \\ x \\ x^2 \end{bmatrix} = a + bx + cx^2$$

Eg. XOR: $x \in (\pm 1, \pm 1)$ $y = x_1 x_2$



feature: $\phi(x) = (1, x_1, x_2, x_1^2, x_2^2, x_1 x_2)$

$$w = \begin{bmatrix} 0 \\ 0 \\ 0 \\ 0 \\ 0 \\ 1 \end{bmatrix}$$

Feature is problem-specific

Feature expansion:

Reduce bias error

$$\phi(x) = \begin{cases} e_i & x_i = x \\ 0 & \text{otherwise} \end{cases}$$

$$y_i \in \{-1, +1\} \quad \text{we} \quad \begin{bmatrix} y_1 \\ \vdots \\ y_n \end{bmatrix}^T \begin{bmatrix} \phi(x_1) \\ \vdots \\ \phi(x_n) \end{bmatrix} = \begin{bmatrix} y_1 x_1 \\ \vdots \\ y_n x_n \\ 0 \end{bmatrix}$$

But feature expansion can do bad on future data.

Eg. $(x, y) \quad x \sim U(0, 1) \quad y \sim N(0, 1)$

$$x \sim \mathcal{N}^T(1, x, x^2, \dots, x^r)$$

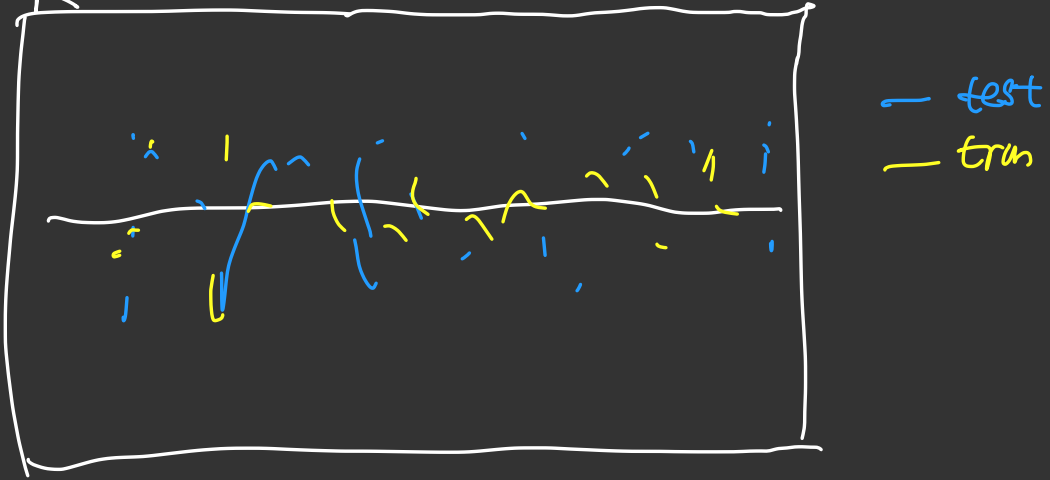
Method: OLS solution \hat{w}_{OLS}

$$\left\| \begin{bmatrix} -\phi(x_1)^T \\ \vdots \\ -\phi(x_n)^T \end{bmatrix} w - \begin{bmatrix} y_1 \\ \vdots \\ y_n \end{bmatrix} \right\|^2 = 0$$

$$\begin{matrix} r \rightarrow r+1 = n & n \leq r+1 \\ \text{as } \# \text{ of points} \end{matrix}$$

\Rightarrow Maybe we can make $r \rightarrow \infty$ so this model can fit to inf points?

$$\hat{R}_1 = 0.069204 \quad R_1 = 0.081745 \quad \hat{R}_2 = 0.006957 \quad R_2 = 0.057077$$



Degree 1: training error 0.069, test error 0.0817

Degree 16: training error 0.0031, test error 0.05706

Overfitting: when training error good,
but bad on testing error.

① often it means a model is complicated in a way that is compatible with the data

② a bias / variance tradeoff

③ train error random test error
not mean "complex models are bad", DV can perform good

Resolution for overfitting

① Reduce model complexity via

regularized ERM: Pick λ so and

$$\min_{w \in \mathcal{H}} \hat{R}(w) + \underbrace{\frac{\lambda}{2} \|w\|_2^2}_{\text{regularization (penalty) (weight decay)}}$$

regularization (penalty) (weight decay)

not very large optimization)

$$w = (X^T X + \underbrace{\lambda n I}_T)^{-1} X^T y$$

hyperparameter

$$\nabla_w \left(\hat{R}(w) + \frac{\lambda}{2} \|w\|_2^2 \right) = \lambda_n \hat{R}(w) + \lambda w$$

Loss functions and multi-class prediction

Standard Classification losses

① hinge : $L_{\text{hinge}}(z) = \max\{0, 1 - z\}$

② squared: $L_{\text{sqr}}(z) = (1 - z)^2$

③ logistic: $L_{\text{logistic}}(z) / \ln(2) = (\ln(1 + e^{-z})) / \ln(2)$

④ exponential: $\exp(z) = e^{-z}$

Design losses with MLE

$$\arg \max_w \prod_{i=1}^n P_w(Y_i | X_i) = \arg \min_w - \ln \prod_{i=1}^n P_w(Y_i | X_i) = \arg \min_w \ln \frac{1}{\prod_{i=1}^n P_w(Y_i | X_i)}$$

Squared loss $P_w(Y|X) = \text{standard Gaussian with mean } X^T w$

$$\ln P_w(Y|X) = \ln \frac{\exp(-\frac{1}{2\sigma^2} (y - w^T X)^2)}{\sigma \sqrt{2\pi}}$$

$$= \frac{-1}{2\sigma^2} (y - w^T X)^2 - \frac{N}{2} \ln(2\pi\sigma^2) \Rightarrow \arg \max_w P_w(Y|X) = \arg \min_w \pm (w^T X)$$

least square

logistic loss

$$p_w(y=1|x) = \frac{1}{1 + \exp(-x^T w)}$$

← sigmoid distribution

$$p_w(y=-1|x) = 1 - p_w(y=1|x) = \frac{\exp(-x^T w)}{1 + \exp(-x^T w)} = \frac{1}{1 + \exp(x^T w)}$$

Then

$$\begin{aligned} \ln \frac{1}{p_w(y|x)} &= \ln \frac{1}{p_w(y=1|x)^{(1+y)/2} (1 - p_w(y=1|x))^{(1-y)/2}} \\ &= \frac{1+y}{2} \ln(1 + \exp(-x^T w)) + \frac{1-y}{2} \ln(1 + \exp(x^T w)) \\ &= \ln(1 + \exp(-yx^T w)) \end{aligned}$$

$$\arg \min_w \ln \prod_{i=1}^n \frac{1}{p_w(y_i|x_i)} = \arg \min_w \frac{1}{n} \sum_{i=1}^n \ln(1 + \exp(-y_i x_i^T w))$$