

CMPUT 267 Basics of Machine Learning

Prediction and Optimal Predictors Linear Regression



February 29, 2024

Announcements

- ▷ Participation and Reading Exercises for MLE and MAP
 - ▷ on eClass tonight along with recorded lecture covering MLE and MAP.
 - ▷ Exercises will be up for a week.
- ▷ Assignment 2a and 2b : deadlines pushed. Check eClass and course website.

Outline

1. Recap
2. Optimal Prediction for Classification - Example
3. Irreducible vs. Reducible Error
4. MLE Formulation for Linear Regression

Recap

- ▶ **Supervised learning problem:** Learn a **predictor** $f : \mathcal{X} \rightarrow \mathcal{Y}$ from a dataset $\mathcal{D} = \{(\mathbf{x}_i, y_i)\}_{i=1}^n$.
 - ▶ \mathcal{X} is the set of **observations** and \mathcal{Y} is the set of **targets**.
- ▶ **Classification** problems have discrete targets.
- ▶ **Regression** problems have continuous targets (order matters).

Recap: Optimal Prediction

- ▷ Suppose we know the true joint distribution $p(\mathbf{x}, y)$ and we want to use it to make predictions in a classification problem.
- ▷ The **optimal classification predictor** makes the best use of this function.
- ▷ As with the optimal estimator, we measure the quality of a predictor $f(\mathbf{x})$ by its **expected cost** $\mathbb{E}[C]$.
- ▷ The optimal predictor **minimizes** $\mathbb{E}[C]$.

$$\mathbb{E}[C] = \int_{\mathcal{X}} \sum_{y \in \mathcal{Y}} \text{cost}(\overset{\hat{y}}{f(\mathbf{x})}, y) p(\mathbf{x}, y) d\mathbf{x},$$

where $\text{cost}(\hat{y}, y)$ is the cost of predicting \hat{y} when the true value is y , and $C = \text{cost}(f(\mathbf{x}), y)$ is a random variable.

Recap: Optimal Classification Prediction

▷ Bayes risk classifier:

$$f^* = \arg \min_{f \in \mathcal{F}} \int_{\mathcal{X}} p(\mathbf{x}) \mathbb{E}[C | \mathbf{X} = \mathbf{x}] d\mathbf{x}$$

$$f^*(\mathbf{x}) = \arg \min_{f \in \mathcal{F}} \mathbb{E}[C | \mathbf{X} = \mathbf{x}] = \arg \min_{\hat{y} \in \mathcal{Y}} \sum_{y \in \mathcal{Y}} \text{cost}(\hat{y}, y) p(y | \mathbf{x})$$

▷ 0 – 1 cost function:

∪

$$f^* = \arg \max_{\hat{y} \in \mathcal{Y}} p(y | \mathbf{x}) \quad \text{mode of } p(y | \mathbf{x})$$

Example

0-1 cost, multiple values for y $y \in \{1, 2, 3, 4\}$

$$\text{cost}(\hat{y}, y) = \begin{cases} 0 & \text{if } \hat{y} = y \\ 1 & \text{else} \end{cases}$$

$$f^*(x) = \underset{\hat{y} \in Y}{\text{argmin}} E[\text{cost} | X=x] = \underset{\hat{y} \in Y}{\text{argmax}} \underbrace{P(y|x)}$$

$$E[\text{cost} | X=x]$$

$$= 0.6 \text{cost}(\hat{y}, 1)$$

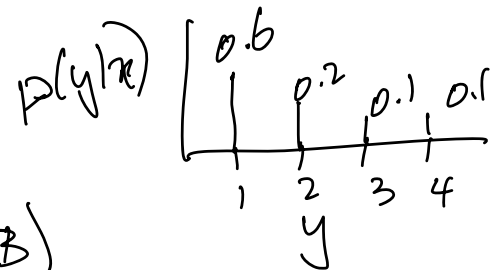
$$+ 0.2 \text{cost}(\hat{y}, 2) + 0.1 \text{cost}(\hat{y}, 3)$$

$$+ 0.1 \text{cost}(\hat{y}, 4)$$

$$\text{if } \hat{y} = 1 : E[\text{cost} | X=x] = 0.6 \cdot 0 + 0.2 \cdot 1 + 0.1 \cdot 1 + 0.1 \cdot 1 = 0.4$$

$$\text{if } \hat{y} = 2$$

$$= 0.6 \cdot 1 + 0.2 \cdot 0 + 0.1 \cdot 1 + 0.1 \cdot 1 = 0.8 \dots$$



$\hat{y} = 1$ is lowest cost

Recap: Optimal Regression Prediction

▷ Most common cost functions:

1. squared error: $\text{cost}(\hat{y}, y) = (\hat{y} - y)^2$

2. absolute error: $\text{cost}(\hat{y}, y) = |\hat{y} - y|$

▷ Squared error function penalizes large error values more than absolute error function.

▷ Optimal prediction for squared error:

$$f^*(\mathbf{x}) = \mathbb{E}[Y \mid \mathbf{X} = \mathbf{x}]$$

For absolute error $f^*(x) = \text{median}(Y|x)$

Example: Classification

cr

- ▷ A medical example where “0 – 1” cost, but with cost differing by type of wrong answer.
- ▷ $y = -1$: no disease; $y = 1$: disease.
- ▷ if disease predicted, further tests; no tests if no disease predicted.
- ▷ **false positive**: leads to unnecessary test.
- ▷ **false negative**: leads to untreated disease (and law suit later).

		y - true	
		-1 (no disease)	1 (disease)
\hat{y} prediction	-1 (no disease)	0 TN	999 FN
	1 (disease)	1 FP	1 TP

0

Example

- ▷ $y \in \{-1, 1\}$
- ▷ Given \mathbf{x} , let:

$$p(y = 1 | \mathbf{x}) = p_1 \quad p(y = -1 | \mathbf{x}) = p_0 \quad (p_0 = 1 - p_1)$$

$$f^*(\mathbf{x}) = \arg \min_{f \in \mathcal{F}} \mathbb{E}[C | \mathbf{X} = \mathbf{x}]$$

$$f^*(x) = \arg \min_{\hat{y} \in \mathcal{Y}} \sum_{y \in \mathcal{Y}} \text{cost}(\hat{y}, y) p(y|x)$$

$$f^* = \underset{\hat{y} \in \mathcal{Y}}{\operatorname{argmin}} \sum_{\substack{y \in \mathcal{Y} \\ \hat{y} = +1 \text{ or } -1}} \operatorname{cost}(\hat{y}, y) p(y|x) = \underset{\hat{y} \in \mathcal{Y}}{\operatorname{argmin}} \left[\operatorname{cost}(\hat{y}, y=1) p(y=1|x) + \operatorname{cost}(\hat{y}, y=-1) p(y=-1|x) \right]$$

$$= \underset{\hat{y} \in \mathcal{Y}}{\operatorname{argmin}} \underbrace{\operatorname{cost}(\hat{y}, y=1) p_1 + \operatorname{cost}(\hat{y}, y=-1) p_0}_{L(\hat{y})}$$

$$\text{if } \hat{y} = -1 \quad L(\hat{y} = -1) = \infty \cdot p_1 + 0 \cdot p_0 = \infty \cdot p_1$$

$$\text{if } \hat{y} = +1 \quad L(\hat{y} = +1) = 1 \cdot p_1 + 1 \cdot p_0 = 1$$

$$\text{if } \infty \cdot p_1 < 1 \quad f^* : \hat{y} = -1 \quad \left[\begin{array}{l} \text{'if } p_1 < \frac{1}{\infty} \\ \infty \end{array} \right]$$

$$\text{else} \quad f^* : \hat{y} = +1 \quad p_1 = P(y=1|x)$$

$$f^* = \underset{\hat{y} \in \mathcal{Y}}{\operatorname{argmin}} L(\hat{y})$$

Error in prediction

Regression

- ▷ Optimal prediction doesn't mean error = 0.
- ▷ What is the quality of our predictor? It may be optimal or suboptimal. Let's look at the expected squared error.
- ▷ First let's consider the optimal predictor, $f^*(\mathbf{x}) = \mathbb{E}[Y | \mathbf{X} = \mathbf{x}]$.

$$\begin{aligned}\mathbb{E}[C] &= \int_{\mathcal{X}} p(\mathbf{x}) \int_{\mathcal{Y}} (f^*(\mathbf{x}) - y)^2 p(y | \mathbf{X} = \mathbf{x}) dy dx \\ &= \int_{\mathcal{X}} p(\mathbf{x}) \int_{\mathcal{Y}} (\mathbb{E}[Y | \mathbf{X} = \mathbf{x}] - y)^2 p(y | \mathbf{X} = \mathbf{x}) dy dx\end{aligned}$$

Error in prediction

- ▷ Optimal prediction doesn't mean error = 0.
- ▷ What is the quality of our predictor? It may be optimal or suboptimal. Let's look at the expected squared error.
- ▷ First let's consider the optimal predictor, $f^*(\mathbf{x}) = \mathbb{E}[Y | \mathbf{X} = \mathbf{x}]$.

$$\begin{aligned}\mathbb{E}[C] &= \int_{\mathcal{X}} p(\mathbf{x}) \int_{\mathcal{Y}} (f^*(\mathbf{x}) - y)^2 p(y | \mathbf{X} = \mathbf{x}) dy dx \\ &= \int_{\mathcal{X}} p(\mathbf{x}) \int_{\mathcal{Y}} (\mathbb{E}[Y | \mathbf{X} = \mathbf{x}] - y)^2 p(y | \mathbf{X} = \mathbf{x}) dy dx\end{aligned}$$

$$= \int_{\mathcal{X}} p(\mathbf{x}) \text{Var}[Y | \mathbf{X} = \mathbf{x}]$$

- ▷ This is **irreducible error**.

Error for any predictor

$$(a+b)^2$$

Now let's consider the expected square error for a suboptimal predictor, $f(\mathbf{x})$.

$$\begin{aligned}\mathbb{E}[C | \mathbf{X}] &= \mathbb{E} \left[(f(\mathbf{x}) - Y)^2 \mid \mathbf{X} = \mathbf{x} \right] = \mathbb{E} \left[\underbrace{(f(\mathbf{x}) - \mathbb{E}[Y | \mathbf{X} = \mathbf{x}])}_a + \underbrace{\mathbb{E}[Y | \mathbf{X} = \mathbf{x}] - Y}_b \right]^2 \mid \mathbf{X} = \mathbf{x} \\ &= \mathbb{E} \left[\underbrace{(f(\mathbf{x}) - \mathbb{E}[Y | \mathbf{X} = \mathbf{x}])^2}_{a^2} + \underbrace{2(f(\mathbf{x}) - \mathbb{E}[Y | \mathbf{X} = \mathbf{x}])}_{2ab} (\mathbb{E}[Y | \mathbf{X} = \mathbf{x}] - Y) \right. \\ &\quad \left. + \underbrace{(\mathbb{E}[Y | \mathbf{X} = \mathbf{x}] - Y)^2}_{b^2} \mid \mathbf{X} = \mathbf{x} \right]\end{aligned}$$

Error for any predictor

Now let's consider the expected square error for a suboptimal predictor, $f(\mathbf{x})$.

$$\begin{aligned}\mathbb{E}[C | \mathbf{X}] &= \mathbb{E} \left[(f(\mathbf{x}) - Y)^2 \mid \mathbf{X} = \mathbf{x} \right] = \mathbb{E} \left[(f(\mathbf{x}) - \mathbb{E}[Y | \mathbf{X} = \mathbf{x}] + \mathbb{E}[Y | \mathbf{X} = \mathbf{x}] - Y)^2 \mid \mathbf{X} = \mathbf{x} \right] \\ &= \mathbb{E} \left[(f(\mathbf{x}) - \mathbb{E}[Y | \mathbf{X} = \mathbf{x}])^2 + 2(f(\mathbf{x}) - \mathbb{E}[Y | \mathbf{X} = \mathbf{x}])(\mathbb{E}[Y | \mathbf{X} = \mathbf{x}] - Y) \right. \\ &\quad \left. + (\mathbb{E}[Y | \mathbf{X} = \mathbf{x}] - Y)^2 \mid \mathbf{X} = \mathbf{x} \right]\end{aligned}$$

Error (cont'd), middle term

$$\begin{aligned} & \mathbb{E}[(f(\mathbf{x}) - \mathbb{E}[Y | X = \mathbf{x}])(\mathbb{E}[Y | X = \mathbf{x}] - Y) | \mathbf{X} = \mathbf{x}] \\ &= (f(\mathbf{x}) - \mathbb{E}[Y | X = \mathbf{x}]) \mathbb{E}[(\mathbb{E}[Y | X = \mathbf{x}] - Y) | X = \mathbf{x}] \\ &= (f(\mathbf{x}) - \mathbb{E}[Y | X = \mathbf{x}]) (\mathbb{E}[Y | X = \mathbf{x}] - \mathbb{E}[Y | X = \mathbf{x}]) \\ &= (f(\mathbf{x}) - \mathbb{E}[Y | X = \mathbf{x}]) 0 \\ &= 0 \end{aligned}$$

Error for any predictor

$$\mathbb{E}[C] = \int_{\mathbf{x} \in \mathcal{X}} p(\mathbf{x}) \mathbb{E}[C | \mathbf{x} = \mathbf{x}] d\mathbf{x}$$

$$\begin{aligned} \mathbb{E}[C] &= \mathbb{E}[\mathbb{E}[C | \mathbf{X} = \mathbf{x}]] && \begin{matrix} a^2 & b^2 \end{matrix} \\ &= \mathbb{E} \left[\mathbb{E} \left[(f(\mathbf{x}) - \mathbb{E}[Y | \mathbf{X} = \mathbf{x}])^2 + (\mathbb{E}[Y | \mathbf{X} = \mathbf{x}] - Y)^2 \mid \mathbf{X} = \mathbf{x} \right] \right] \\ &= \mathbb{E} \left[(f(\mathbf{X}) - \mathbb{E}[Y | \mathbf{X} = \mathbf{x}])^2 \right] + \mathbb{E} \left[(\mathbb{E}[Y | \mathbf{X} = \mathbf{x}] - Y)^2 \mid \mathbf{X} = \mathbf{x} \right] \\ &= \mathbb{E} \left[(f(\mathbf{X}) - f^*(\mathbf{X}))^2 \right] + \mathbb{E} \left[(f^*(\mathbf{X}) - Y)^2 \right] \end{aligned}$$

Error for any predictor

$$\begin{aligned}\mathbb{E}[C] &= \mathbb{E}[\mathbb{E}[C \mid \mathbf{X} = \mathbf{x}]] \\ &= \mathbb{E} \left[\mathbb{E} \left[(f(\mathbf{x}) - \mathbb{E}[Y \mid \mathbf{X} = \mathbf{x}])^2 + (\mathbb{E}[Y \mid \mathbf{X} = \mathbf{x}] - Y)^2 \mid \mathbf{X} = \mathbf{x} \right] \right] \\ &= \mathbb{E} \left[(f(\mathbf{X}) - \mathbb{E}[Y \mid \mathbf{X} = \mathbf{x}])^2 \right] + \mathbb{E} \left[(\mathbb{E}[Y \mid \mathbf{X} = \mathbf{x}] - Y)^2 \mid \mathbf{X} = \mathbf{x} \right] \\ &= \mathbb{E} \left[(f(\mathbf{X}) - f^*(\mathbf{X}))^2 \right] + \mathbb{E} \left[(f^*(\mathbf{X}) - Y)^2 \right]\end{aligned}$$

reducible error

irreducible error

How to reduce the reducible error?

- ▶ We want to make the error between the f that we learn and the optimal f^* smaller.
- ▶ Let's assumed the hypothesis space we're looking in, \mathcal{F} , is the space of linear functions.
- ▶ **Sources of reducible error**
 1. **Limited hypothesis space.** We assumed linear functions, but maybe f^* is non linear.
 2. **Insufficient optimization.** We might have used gradient descent, but didn't fully optimize f - stopped too early?
 3. **Limited data.** Not enough samples to identify a good f .

How to reduce the reducible error?

1. **Limited hypothesis space.** We assumed linear functions, but maybe f^* is non linear.
 - ▷ Solution: make the hypothesis space bigger (e.g. polynomials?)
2. **Insufficient optimization.** We might have used gradient descent, but didn't fully optimize f - stopped too early?
 - ▷ Solution: set step size and number of epochs more carefully to ensure you're at a stationary point.
3. **Limited data.** Not enough samples to identify a good f .
 - ▷ Solution: gather more data.

How to reduce the irreducible error?

$$\int p(x) \text{Var}[Y|X=x]$$

- ▷ It's **irreducible**...
- ▷ It's the variance of Y given X : $\text{Var}(Y | \mathbf{X} = \mathbf{x})$.
- ▷ Improving the learned function **cannot** change the inherent variance in Y .
- ▷ BUT: what is the source of variance in Y given \mathbf{x} ?

How to reduce the irreducible error?

- ▷ It's **irreducible**...
- ▷ It's the variance of Y given X : $\text{Var}(Y | \mathbf{X} = \mathbf{x})$.
- ▷ Improving the learned function **cannot** change the inherent variance in Y .
- ▷ BUT: what is the source of variance in Y given \mathbf{x} ?
 - ▷ partial observability
 - ▷ stochasticity in the system

Linear Predictors

Linear Regression

Setting: $D = \{(\mathbf{x}_i, y_i)\}_{i=1}^n$ $\mathbf{x}_i \in \mathbb{R}^d$ $y_i \in \mathbb{R}$

$$f(\mathbf{x}_i) = \underline{w_0} + w_1 x_{i,1} + w_2 x_{i,2} + \dots + w_d x_{i,d}$$

$$f(\mathbf{x}_i) = \sum_{j=0}^d w_j x_{i,j} \quad (x_{i,0} = 1)$$

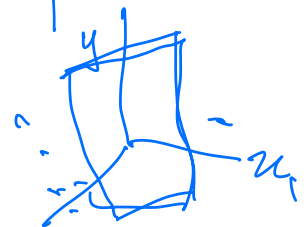
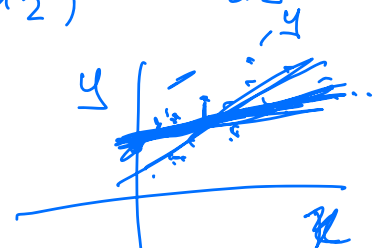
$P(y|\mathbf{x})$

$d=1 \quad w_0=0$

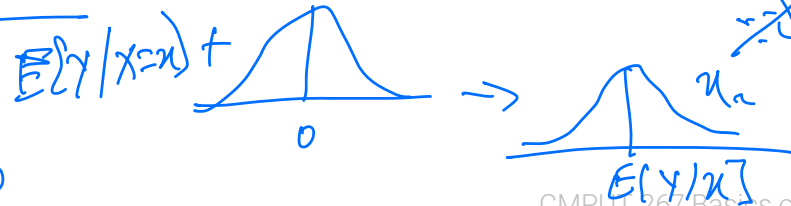
$$P(y|\mathbf{x}) = N(\underline{\mathbf{x}\mathbf{w}}, \sigma^2)$$

$$\underline{P(y|\mathbf{x})} \sim N(\mu = f(\mathbf{x}), \sigma^2)$$

relative $\bar{\mathbf{x}} : [1, x_1, x_2, \dots, x_d]$



$$y_i = \underbrace{\sum_{j=0}^d w_j x_{i,j}}_{\mathbf{w}^T \mathbf{x}} + \underline{\varepsilon_i} \quad \varepsilon_i \sim N(0, \sigma^2)$$



$$\mathbf{w}^T \mathbf{x} = \bar{\mathbf{x}}^T \mathbf{w}$$

Hypothesis space:

$$\mathcal{F} = \left\{ p(\cdot | x) = \mathcal{N}(w^T x, \sigma^2) \mid w \in \mathbb{R}^d \right\}$$

Goal: find parameters $w_j, j=0, \dots, d$
 $\vec{w} \in \mathbb{R}^{d+1}$

that is optimal in \mathcal{F}

$$\underline{\text{MLE}} : \underset{w \in \mathbb{R}^{d+1}}{\text{argmin}} \quad \frac{1}{n} \sum_{i=1}^n \underline{-\ln p(y_i | x_i, \vec{w})}$$

MLE for Linear regression

$$\operatorname{argmax}_{w \in \mathbb{R}^{d+1}} p(D|w) = \operatorname{argmin}_{w \in \mathbb{R}^{d+1}} -\ln(p(D|w))$$

$$l(w) = \frac{1}{n} \sum_{i=1}^n l_i(w)$$

$$= \operatorname{argmin}_{w \in \mathbb{R}^{d+1}} \frac{1}{n} \left(\sum_{i=1}^n -\ln P(y_i | \tilde{x}_i, \tilde{w}) \right) \quad N(\tilde{x}^T w, \sigma^2)$$

$$= \operatorname{argmin}_{w \in \mathbb{R}^{d+1}} -\frac{1}{n} \sum_{i=1}^n \ln \frac{1}{\sqrt{2\pi\sigma^2}} \exp\left(-\frac{(y_i - \tilde{x}^T w)^2}{2\sigma^2}\right)$$

$$= \operatorname{argmin}_{w \in \mathbb{R}^{d+1}} \frac{1}{n} \sum_{i=1}^n \frac{(\tilde{x}^T w - y_i)^2}{2} = \operatorname{argmin}_{w \in \mathbb{R}^{d+1}} \frac{1}{n} \sum_{i=1}^n \frac{1}{2} (\tilde{x}^T w - y_i)^2$$

Ordinary Least Squares

Find ω_{MLE}

$$\triangleright C(\omega) = 0$$

$$\begin{bmatrix} \frac{\partial C(\omega)}{\partial \omega_0} \\ \vdots \\ \frac{\partial C(\omega)}{\partial \omega_d} \end{bmatrix}$$

$$\frac{\partial C(\omega)}{\partial \omega_j} = \frac{1}{n} \sum \frac{\partial C(\omega)}{\partial \omega_j}$$

$$\hookrightarrow \frac{\partial C(\omega)}{\partial \omega_j} = \frac{\partial}{\partial \omega_j} \left(\frac{1}{2} (\mathbf{x}^T \omega - y_i)^2 \right) \frac{\partial (\mathbf{x}^T \omega)}{\partial \omega_j}$$

$$\frac{1}{n} \sum_{i=1}^n (\mathbf{x}^T \omega - y_i) \mathbf{x}_i = 0$$

$$\hookrightarrow \frac{\partial (\mathbf{x}^T \omega)}{\partial \omega_j} = \frac{\partial}{\partial \omega_j} \sum_{k=0}^d \kappa_{ik} \omega_k$$

$$= \sum_{k=0}^d \kappa_{ik} \frac{\partial \omega_k}{\partial \omega_j} = \kappa_{ij}$$

$$\frac{1}{n} \sum (x_i^T \omega - y_i) x_{i0} = 0$$

$$\frac{1}{n} \sum (x_i^T \omega - y_i) x_{i1} = 0$$

⋮

$$\frac{1}{n} \sum (x_i^T \omega - y_i) x_{id} = 0$$

$$\rightarrow \underline{A\omega = b}$$

$$A = \frac{1}{n} \sum_{i=1}^n x_i x_i^T$$

$$b = \frac{1}{n} \sum_{i=1}^n x_i y_i$$

$O(nd^2)$ matrix inversion $O(d^3)$

$$O(nd^2 + d^3)$$

Numerical Solution

SLD.

1. Initialize \vec{w}_0

\vec{w}_t

$$w_{0j} = N(0, 1)$$

2. Sample minibatch (shuffle)

b: batch size

$$3. \vec{w}_{t+1} = \vec{w}_t - \eta_t \underbrace{\frac{1}{b} \sum_{i=1}^b (x_i^T \vec{w}_t - y_i) x_i}_{g_t}$$

η_t (vector)

numerical issues at initialization

$$\eta_{t,j} = \frac{1}{10^4 \sqrt{g_{t,j}}}$$

w_{MLE}

$$\sim 10^4 t$$

$$\bar{g}_{t,j} = \bar{g}_{t-1,j} + g_{t,j}^2$$

$$\bar{g}_0 = 0$$

$O(kbd)$
(for k iterations of SLD)

$$\underline{\underline{\omega_{t+1}}} = \omega_t - \eta_t \underbrace{c'(\omega_t)}_{\left. \frac{1}{n} \sum_{i=1}^n c'_i(\omega_t) \right\}}$$

$b \ll n$
 $\frac{1}{b} \sum_{i=1}^b c'_i(\omega)$