

# Maximum Likelihood, MMSE



- maximum likelihood estimation
  - regression
  - MMSE

# Big Picture

- machine learning is about **learning functions from data**

$$\hat{y} = f(\mathbf{x})$$

- the inputs to the functions are called **features**
- the true output is called a label

$$\mathbf{x} = \boxed{9}$$

classification:  $y \in \{0, 1, \dots, 9\}$

regression:  $y \in \mathbb{R}$

$$\mathcal{D} = \{(\mathbf{x}_i, y_i)\}_{i=1}^n \quad (\mathbf{x}_i, y_i) \stackrel{i.i.d.}{\sim} p(\mathbf{x}, y)$$

- if we have a good approximation of  $p(\mathbf{x}, y)$ , we know what to do:

$$\hat{y} = \arg \max_y p(y|\mathbf{x}) \quad \text{MAP estimate}$$

$$\hat{y} = \arg \max_y p(\mathbf{x}|y)p(y) \quad \text{posterior } \propto \text{likelihood } \times \text{ prior}$$

$$\hat{y} = \arg \max_y p(\mathbf{x}|y) \quad \text{ML (maximum likelihood) estimate}$$

$$\hat{y} = \arg \max_y p_y(\mathbf{x})$$

# Maximum Likelihood Example

- flip a coin with unknown bias  $n$  times       $H, H, H, T$
- how can we estimate bias of the coin  $\theta$ ?

$$X_i \stackrel{i.i.d.}{\sim} \text{Bernoulli}(\theta)$$

$$p(x) = \begin{cases} \theta & x = 1 \\ (1 - \theta) & x = 0 \end{cases}$$

$$p(\mathbf{x}) = p(x_1 = 1)p(x_2 = 1)p(x_3 = 1)p(x_4 = 0)$$

$$p_\theta(\mathbf{x}) = \theta^3(1 - \theta)$$

- think of  $p_\theta(\mathbf{x})$  as a family of distributions parameterized by  $\theta$

$$p_\theta(\mathbf{x}) = \prod_{i=1}^n p_\theta(x_i)$$

$$p_\theta(\mathbf{x}) = \theta^k(1 - \theta)^{n-k}$$

where  $k$  is number of heads

- answer should be  $\hat{\theta} = \frac{k}{n}$

# Maximum Likelihood

- observe some data  $\mathbf{x}$

$$\hat{y} = \arg \max_y p(\mathbf{x}|y)p(y)$$

- family of distributions parametrized by  $\theta$

$$p_\theta(\mathbf{x})$$

$$\hat{y} = \arg \max_y p(\mathbf{x}|y)$$

$$\hat{\theta} = \arg \max_\theta p_\theta(\mathbf{x})$$

$$\hat{y} = \arg \max_y p_y(\mathbf{x})$$

$$\hat{\theta} = \arg \max_\theta p_\theta(\mathbf{x})$$

- $L(\theta) = p_\theta(\mathbf{x})$  is the likelihood function (fixed  $\mathbf{x}$ , function of  $\theta$ )
- $\ell(\theta) = \log(p_\theta(\mathbf{x}))$  is the log-likelihood function (fixed  $\mathbf{x}$ , function of  $\theta$ )

$$\hat{\theta} = \arg \max_\theta L(\theta) = \arg \max_\theta \ell(\theta)$$

$\hat{\theta}$  is the maximum likelihood estimate of  $\theta$

# Maximum Likelihood for Mean of Gaussian

- model:  $X_1, X_2, \dots, X_n \stackrel{i.i.d.}{\sim} \mathcal{N}(\theta, 1)$  for unknown  $\theta$

$$L(\theta) = p_\theta(\mathbf{x}) = \prod_{i=1}^n p_\theta(x_i) = \prod_i \frac{1}{\sqrt{2\pi}} e^{-\frac{(x_i - \theta)^2}{2}}$$

‘likelihood of  $\theta$  given we saw the data  $\mathbf{x}$ ’

‘what’s the most likely  $\theta$  given we saw the data  $\mathbf{x}$ ?’

$$\arg \max_{\theta} L(\theta) = \arg \max_{\theta} \prod_i \frac{1}{\sqrt{2\pi}} e^{-\frac{(x_i - \theta)^2}{2}}$$

```
x = np.random.randn(100, 1) + mu
print(x)
```

$\mathbf{x} =$   
[ [-0.10223044]  
[ 0.74557181]  
[ 1.61851936]  
[-1.05609849]  
[-0.50006279]  
[-1.24555422]  
...]

$$= \arg \min_{\theta} \sum_i (x_i - \theta)^2$$

$$\hat{\theta}_{\text{ML}} = \frac{1}{n} \sum_i x_i$$

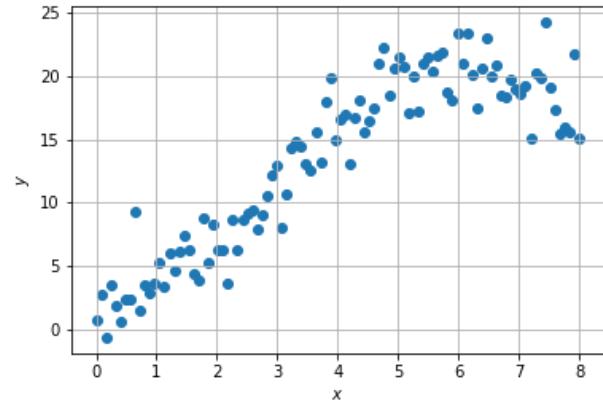
# Regression and Maximum Likelihood

- model  $y$  as a weighted sum of features  $\mathbf{x}$ , plus noise  $Z$

$$y_i = \mathbf{w}^T \mathbf{x}_i + Z_i \quad Z_i \sim \mathcal{N}(0, \sigma^2)$$

$$y_i \sim \mathcal{N}(\mathbf{w}^T \mathbf{x}_i, \sigma^2)$$

$$\mathcal{D} = \{(\mathbf{x}_i, y_i)\}_{i=1}^n$$



- what's the likelihood of  $\mathbf{w}$  given we saw  $\mathcal{D}$ ?

$$L(\mathbf{w}) = p_{\mathbf{w}}(\mathbf{y}) = \prod_{i=1}^n \frac{1}{\sqrt{2\pi\sigma^2}} e^{-\frac{(y_i - \mathbf{x}_i^T \mathbf{w})^2}{2\sigma^2}}$$

- goal: find coefficients  $\mathbf{w}$  so that

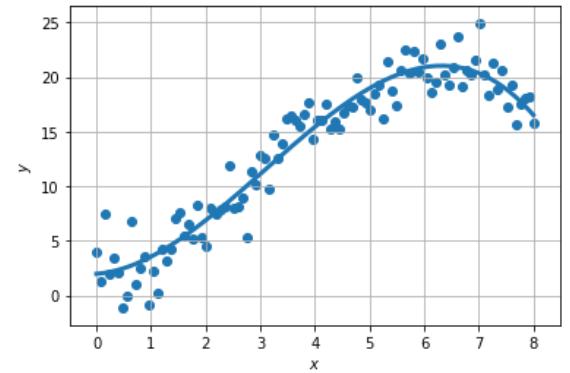
$$y = w_1 x^3 + w_2 x^2 + w_3 x^1 + w_4$$

$$y_i = \mathbf{w}^T \mathbf{x}_i$$

$$\mathbf{x}_i^T = [x_i^3 \ x_i^2 \ x_i^1 \ 1]$$

# Regression and Maximum Likelihood

$$L(\mathbf{w}) = p_{\mathbf{w}}(\mathbf{y}) = \prod_{i=1}^n \frac{1}{\sqrt{2\pi\sigma^2}} e^{-\frac{(y_i - \mathbf{x}_i^T \mathbf{w})^2}{2\sigma^2}}$$
$$\mathbf{w}_{\text{ML}} = \arg \max_{\mathbf{w}} L(\mathbf{w})$$



$$= \arg \min_{\mathbf{w}} \sum_i (y_i - \mathbf{w}^T \mathbf{x}_i)^2$$

$$\mathbf{x}^T = [x^3 \ x^2 \ x^1 \ 1]$$

$$= \arg \min_{\mathbf{w}} \|\mathbf{X}\mathbf{w} - \mathbf{y}\|_2^2$$

$$\mathbf{w}_{\text{ML}} = (\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T \mathbf{y}$$

- maximum likelihood with gaussian noise gives rise to the squared error loss

# Minimum Mean Square Error Estimation

# Minimum Mean Squared Error

recap . . .

- MAP estimate:  $\hat{y} = \arg \max_y p(\mathbf{x}|y)p(y)$
- ML estimate:  $\hat{y} = \arg \max_y p_y(\mathbf{x})$ 
  - **claim:** MAP minimizes risk under 0/1 loss.
  - **claim:** ML minimizes risk under 0/1 loss and uniform prior.

$$E[\ell(\hat{y}, y)] = E[\mathbb{1}_{\{\hat{y} \neq y\}}] = \mathbb{P}(\hat{y} \neq y)$$

MAP by definition picks the most probable  $y$ .

- MMSE (Minimum mean squared error) estimate:

$$E[\ell(\hat{y}, y)] = E[(\hat{y} - y)^2]$$

- **claim (obvious):** MMSE minimizes risk under squared error loss.

# Minimum Mean Squared Error Estimation

- we observe features  $\mathbf{x}$  and want to estimate  $y$ , and we know  $p(\mathbf{x}, y)$
- find an estimate  $\hat{y}$  that minimizes risk with squared error loss:

$$E[\ell(y, \hat{y})] = E[(y - \hat{y})^2]$$

- claim: the MMSE estimate of  $y$  given the features  $\mathbf{x}$  and complete knowledge of  $p(\mathbf{x}, y)$  is:

$$\hat{y} = E[y|\mathbf{x}]$$

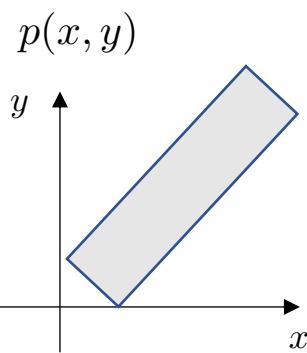
- note: MMSE is basis for denoising techniques such as Wiener filter, Kalman Filter.

# MMSE Example

- claim: the MMSE estimate of  $y$  given the features  $\mathbf{x}$  and complete knowledge of  $p(\mathbf{x}, y)$  is:

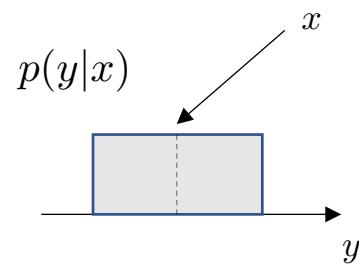
$$\hat{y} = f(\mathbf{x}) = E[y|\mathbf{x}]$$

- example:



- given  $x$ , what is your best guess for  $y$ ?

$$y|x \sim U[x - \Delta, x + \Delta]$$



- MMSE:  $\hat{y} = E[y|x] = x$
- MAP:  $\hat{y}$  can be anything between  $x - \Delta$  and  $x + \Delta$

# MMSE Proof

- claim: the MMSE estimate of  $y$  given the features  $\mathbf{x}$  and complete knowledge of  $p(\mathbf{x}, y)$  is:

$$\hat{y} = f(\mathbf{x}) = E[y|\mathbf{x}]$$

- proof:

1. note that  $\min_a E[(y - a)^2]$  is given by  $a = E[y]$  since

$$\begin{aligned} E[(y - a)^2] &= E[((y - E[y]) + (E[y] - a))^2] \\ &= E[((y - E[y])^2 + (E[y] - a)^2 + 2(y - E[y])(E[y] - a))] \\ &= E[(y - E[y])^2] + (E[y] - a)^2 \\ &\geq E[(y - E[y])^2] \quad \text{with equality if and only if } E[y] = a \end{aligned}$$

2. for fixed  $\mathbf{x}$ , we have that  $E_y[(y - f(\mathbf{x}))^2 | \mathbf{x}]$  is minimized by  $f(\mathbf{x}) = E[y|\mathbf{x}]$

$$E[(y - f(\mathbf{x}))^2] = E_{\mathbf{x}}[E_y[(y - f(\mathbf{x}))^2 | \mathbf{x}]]$$

$f(\mathbf{x}) = E[y|\mathbf{x}]$  minimizes mean squared error.

# MMSE

- claim: the minimum risk under the squared error loss given the features  $\mathbf{x}$

and complete knowledge of  $p(\mathbf{x}, y)$  is:

$$E_{\mathbf{x}}[\text{var}(y|\mathbf{x})]$$

- the MMSE (the error, not the estimate) of  $y$  given the features  $\mathbf{x}$  and complete knowledge of  $p(\mathbf{x}, y)$  is:

$$E_{\mathbf{x}}[\text{var}(y|\mathbf{x})]$$

- proof: