

Lecture 6

Regularization

MLE can overfit the training data.

Ex. $\vec{\theta} = (\theta_h, 1 - \theta_h)$
↑ heads in coin tosses

$$\mathcal{D} = \underbrace{\{h, h, h\}}_{N=3} \Rightarrow \theta_h^* = \frac{3}{3} = 1$$

Now, $\text{Ber}(y | \vec{\theta}^*)$ will predict no tails \Rightarrow the model does not generalize

To prevent overfitting, one uses regularization:

$$J(\vec{\theta}; \lambda) = + \sum_{n=1}^N \log p(\tilde{y}_n | \tilde{x}_n, \vec{\theta}) + \lambda \underbrace{C(\vec{\theta})}_{\text{complexity penalty}}$$

If $\lambda = 1$ and $C(\vec{\theta}) = \log p(\vec{\theta})$, we have prior

$$J(\vec{\theta}; \lambda) = \sum_{n=1}^N \log p(\tilde{y}_n | \tilde{x}_n, \vec{\theta}) + \log p(\vec{\theta}) = \\ = \log p(\mathcal{D} | \vec{\theta}) + \log p(\vec{\theta}) = \log p(\vec{\theta} | \mathcal{D}) + \text{const}$$

MAP estimate: $\vec{\theta}_{\text{MAP}} = \underset{\vec{\theta}}{\operatorname{argmax}} \log p(\vec{\theta} | \mathcal{D})$

Ex. Coin tossing:

use $p(\theta) = \text{Beta}(\theta | a, b)$, where

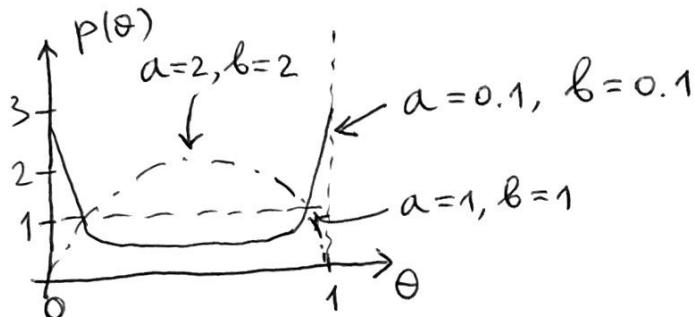
$\text{Beta}(\theta | a, b) = \frac{1}{B(a, b)} \theta^{a-1} (1-\theta)^{b-1}$ is
the beta distribution.

Here, $B(a, b) = \frac{\Gamma(a)\Gamma(b)}{\Gamma(a+b)}$ is the
beta function, and

$$\Gamma(a) = \int_0^\infty dx x^{a-1} e^{-x}.$$

↖ gamma function

Note that $a, b > 0$; if $a=b=1$, the beta distribution is uniform.



$$\underbrace{E[\theta]}_{\text{mean}} = \frac{a}{a+b}$$

We need to maximize

$$\log p(D|\vec{\theta}) + \log p(\vec{\theta}) \Rightarrow \text{maximize}$$

$$\Rightarrow \underbrace{N_1}_{\# \text{ heads}} \log \theta + \underbrace{N_0}_{\# \text{ tails}} \log(1-\theta) + (a-1) \log \theta + (b-1) \log(1-\theta).$$

$$N = N_0 + N_1$$

$$\frac{\partial}{\partial \theta} \log p(\theta | D) \Big|_{\theta_{MAP}} = \frac{N_1 + (a-1)}{\theta_{MAP}} - \frac{N_0 + (b-1)}{1-\theta_{MAP}} = 0, \text{ or}$$

$$[N_1 + (a-1)](1-\theta_{MAP}) = [N_0 + (b-1)]\theta_{MAP},$$

$$\theta_{MAP} = \frac{N_1 + (a-1)}{N + a+b-2}.$$

we might set $a=b=2 \Rightarrow \theta_{MAP} = \frac{N_1 + 1}{N + 2}$
 "add one smoothing"

Ex. Regression:

use zero-mean Gaussian for $p(\vec{\omega})$:

maximize $\underbrace{\log p(D|\vec{\omega})}_{\begin{array}{l} \text{precision prm;} \\ \frac{1}{2} \end{array}} - \lambda |\vec{\omega}|^2$ w.r.t. $\vec{\omega}$

$$- \frac{1}{2\sigma^2} \sum_{i=1}^N (y_i - f(\vec{x}_i; \vec{\omega}, b))^2 + \text{const}(\vec{\omega})$$

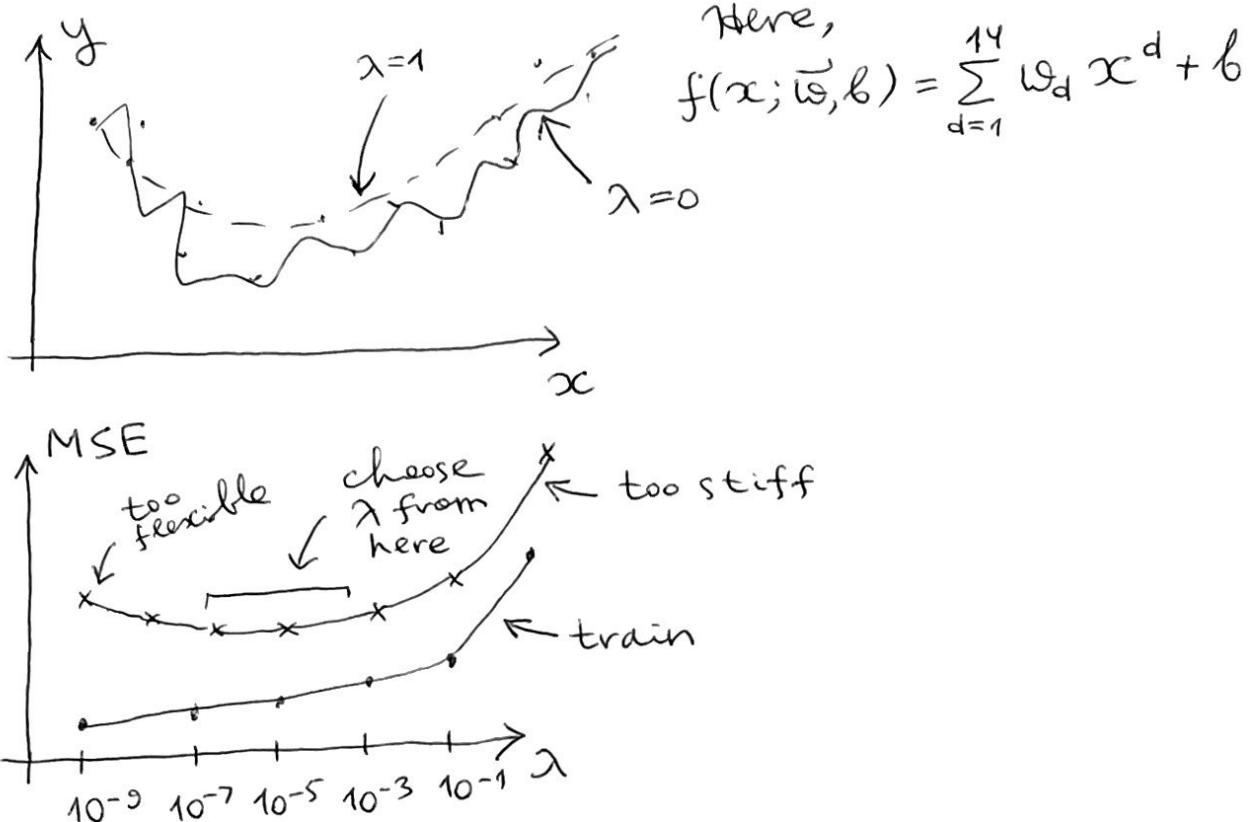
$$\log p(\vec{\omega}) = -\lambda |\vec{\omega}|^2 + \text{const}(\vec{\omega})$$

$$|\vec{\omega}|^2 = \sum_{d=1}^D \omega_d^2 \quad \vec{x} \in \mathbb{R}^D$$

Linear regression:

$$f(\vec{x}; \vec{\omega}, b) = \sum_{d=1}^D \omega_d x_d + b.$$

Choice of λ :



Bayesian statistics

$\frac{\text{likelihood} \ p(\vec{\theta})}{\text{prior}}$

$$p(\vec{\theta} | \mathcal{D}) = \frac{p(\mathcal{D} | \vec{\theta}) p(\vec{\theta})}{\underbrace{\int d\vec{\theta}' p(\mathcal{D} | \vec{\theta}') p(\vec{\theta}')}_{p(\mathcal{D})}, \text{ marginal likelihood}}$$

↑
posterior

Supervised: $\mathcal{D} = \{(\vec{x}_n, \vec{y}_n)\}_{n=1}^N$

Unsupervised: $\mathcal{D} = \{\vec{y}_n\}_{n=1}^N$

Once we have $p(\vec{\theta} | \mathcal{D})$, we can compute the predictive distribution

$$p(\vec{y} | \vec{x}, \mathcal{D}) = \int d\vec{\theta} \ p(\vec{y} | \vec{x}, \vec{\theta}) \ p(\vec{\theta} | \mathcal{D}).$$

↑
supervised case

For example, if $p(\vec{\theta} | \mathcal{D}) = \delta(\vec{\theta} - \vec{\theta}_{MAP})$,

$$p(\vec{y} | \vec{x}, \mathcal{D}) = \underbrace{p(\vec{y} | \vec{x}, \vec{\theta}_{MAP})}_{\text{a single model predominates}}$$

Conjugate priors: priors that are parameterized similarly to the likelihood.

$$\text{Ex. } \begin{cases} p(\mathcal{D} | \theta) = \theta^{N_1} (1-\theta)^{N_0} & \text{likelihood} \\ p(\theta) = \text{Beta}(\theta | a, b) = \frac{1}{B(a,b)} \theta^{a-1} (1-\theta)^{b-1} & \text{prior} \end{cases}$$

$$\text{Then } p(\theta | \mathcal{D}) = \text{Beta}(\theta | a + N_1, b + N_0).$$

In particular, $E[\theta] = \frac{a + N_1}{\underbrace{a+b}_{\tilde{N}} + \underbrace{N_1+N_0}_N}$

\tilde{N} = equire.
sample size

N = sample size

One can write $m = \frac{a}{\tilde{N}}$ θ^*, MLE

$$E[\theta] = \frac{\tilde{N} \frac{a}{\tilde{N}} + N \frac{N_1}{N}}{N + \tilde{N}} = \underbrace{\frac{\tilde{N}}{N + \tilde{N}}}_{\lambda} m + \frac{N}{N + \tilde{N}} \underbrace{\frac{N_1}{N}}_{1-\lambda} = \lambda m + (1-\lambda) \theta^*$$

as $N \uparrow$, $\lambda \downarrow$ and $\underbrace{E[\theta]}_{\text{posterior mean}} \rightarrow \theta^*$.

Predictive distribution:

$$p(y=1 | D) = \int_0^1 d\theta \underbrace{p(y=1|\theta)}_{\theta} p(\theta | D) = E[\theta | D] = \\ = \frac{N_1 + a}{N + a + b}.$$

If $a=b=1$ (uniform prior), we have:

$$p(y=1 | D) = \frac{N_1 + 1}{N + 2}$$

↑ pseudocounts