

## LECTURE 15 - MAXIMUM LIKELIHOOD AND REGULARIZATION

### 1. REGULARIZATION

- **Overfitting/Overtraining:** We discussed overfitting and underfitting. Suppose you have data that you fit a model to, how do you know if you have overfit and/or underfit? *What is Cross Validation?*
- Previously we mentioned two common approaches to avoid overfitting:
  - (1) **More data:** As you have more and more data, it becomes more and more difficult to “memorize” the data and its noise. Often, more data translates to the ability to use a more complex model and avoid overfitting. However, generally, you need exponentially more data with increases to model complexity. So, there is a limit to how much this helps. If you have a very complex model, you need a huge training data set.
  - (2) **Regularization:** Regularization methods add a penalty term to the error function to discourage overfitting. These penalty terms encourage small values limiting the ability to overfit. These penalty terms are a way to trade-off between error and complexity.

$$\begin{aligned}(1) \quad E^*(\mathbf{w}) &= \frac{1}{2} \sum_{n=1}^N (y(x_n, \mathbf{w}) - t_n)^2 + \frac{\lambda}{2} \|\mathbf{w}\|_2^2 \\(2) \quad &= \frac{1}{2} (\mathbf{w}^T \mathbf{X}^T - \mathbf{t}^T) (\mathbf{w}^T \mathbf{X}^T - \mathbf{t}^T)^T + \frac{\lambda}{2} \mathbf{w}^T \mathbf{w}\end{aligned}$$

- *What does each term mean/promote in the minimization? Why does the second term make sense for minimizing complexity?*
- *What happens to  $\mathbf{w}$  with increasing model complexity (and no regularization)?*

$$\begin{aligned}(3) \quad \frac{\partial E^*(\mathbf{w})}{\partial \mathbf{w}} &= 0 = \mathbf{X}^T (\mathbf{w}^T \mathbf{X}^T - \mathbf{t}^T)^T + \frac{\lambda}{2} 2\mathbf{w} \\(4) \quad 0 &= \mathbf{X}^T \mathbf{X} \mathbf{w} - \mathbf{X}^T \mathbf{t} + \lambda \mathbf{w} \\(5) \quad \mathbf{X}^T \mathbf{t} &= (\mathbf{X}^T \mathbf{X} + \lambda \mathbf{I}) \mathbf{w} \\(6) \quad \mathbf{w} &= (\mathbf{X}^T \mathbf{X} + \lambda \mathbf{I})^{-1} \mathbf{X}^T \mathbf{t}\end{aligned}$$

- The  $l_2$  norm penalty is common (one reason it is common: because it works so well mathematically with the least-squares error objective) and, so, has many names: shrinkage, ridge regression, weight decay

- So, what happens when  $\lambda$  is increased? decreased? Can you think of a way to set  $\lambda$ ?
- We looked at the regularization term as a *penalty* term in the objective function. There is another way to interpret the regularization term as well. Specifically, there is a *Bayesian* interpretation.

$$(7) \quad \min E^*(\mathbf{w}) = \max -E^*(\mathbf{w})$$

$$(8) \quad = \max \exp \{-E^*(\mathbf{w})\}$$

$$(9) \quad = \max \exp \left\{ -\frac{1}{2} \sum_{n=1}^N (y(x_n, \mathbf{w}) - t_n)^2 - \frac{\lambda}{2} \|\mathbf{w}\|_2^2 \right\}$$

$$(10) \quad = \max \exp \left\{ -\frac{1}{2} \sum_{n=1}^N (y(x_n, \mathbf{w}) - t_n)^2 \right\} \exp \left\{ -\frac{1}{2} \lambda \|\mathbf{w}\|_2^2 \right\}$$

$$(11) \quad = \max \prod_{n=1}^N \exp \left\{ -\frac{1}{2} (y(x_n, \mathbf{w}) - t_n)^2 \right\} \exp \left\{ -\frac{1}{2} \lambda \|\mathbf{w}\|_2^2 \right\}$$

- So, this is a maximization of the *data likelihood* with a *prior*:  $p(\mathbf{X}|\mathbf{w})p(\mathbf{w})$

## 2. METHOD OF MAXIMUM LIKELIHOOD

- A *data likelihood* is how likely the data is given the parameter set
- So, if we want to maximize how likely the data is to have come from the model we fit, we should find the parameters that maximize the likelihood
- A common trick of maximizing the likelihood is to maximize the log likelihood. Often makes the math much easier. *Why can we maximize the log likelihood instead of the likelihood and still get the same answer?*
- Consider:  $\max \ln \exp \left\{ -\frac{1}{2} (y(x_n, \mathbf{w}) - t_n)^2 \right\}$  We go back to our original objective.

## 3. METHOD OF MAXIMUM A POSTERIORI (MAP)

- Bayes Rule:  $p(Y|X) = \frac{p(X|Y)p(Y)}{p(X)}$
- Consider:  $p(\mathbf{w}|\mathcal{D}) = \frac{p(\mathcal{D}|\mathbf{w})p(\mathbf{w})}{p(\mathcal{D})}$ , i.e., posterior  $\propto$  likelihood  $\times$  prior

## 4. GAUSSIAN DISTRIBUTION

- These are Gaussian distributions:

$$(12) \quad \mathcal{N}(x|\mu, \sigma^2) = \frac{1}{\sqrt{2\pi\sigma^2}} \exp \left\{ -\frac{1}{2} \frac{(x - \mu)^2}{\sigma^2} \right\}$$

- $\sigma^2$  is the variance OR  $\frac{1}{\sigma^2}$  is the *precision*
- So, as  $\lambda$  gets big, variance gets smaller/tighter. As  $\lambda$  gets small, variance gets larger/wider.

- What is the expected value of  $x$ ?

$$(13) \quad E[x] = \int xp(x)dx$$

$$(14) \quad = \int x \frac{1}{\sqrt{2\pi}\sigma^2} \exp\left\{-\frac{1}{2} \frac{(x-\mu)^2}{\sigma^2}\right\} dx$$

- Change of variables: Let

$$(15) \quad y = \frac{x-\mu}{\sigma} \rightarrow x = \sigma y + \mu$$

$$(16) \quad dy = \frac{1}{\sigma} dx \rightarrow dx = \sigma dy$$

- Plugging this into the expectation:

$$(17) \quad E[x] = \int (\sigma y + \mu) \frac{1}{\sqrt{2\pi}\sigma} \exp\left\{-\frac{1}{2}y^2\right\} \sigma dy$$

$$(18) \quad = \int \frac{\sigma y}{\sqrt{2\pi}} \exp\left\{-\frac{1}{2}y^2\right\} dy + \int \frac{\mu}{\sqrt{2\pi}} \exp\left\{-\frac{1}{2}y^2\right\} dy$$

- The first term is an odd function:  $f(-y) = -f(y)$  So,  $E[x] = 0 + \mu = \mu$

## 5. METHOD OF MAXIMUM A POSTERIORI (MAP)

- Bayes Rule:  $p(Y|X) = \frac{p(X|Y)p(Y)}{p(X)}$
- Lets look at this in terms of binary variables, e.g., Flipping a coin:  $X = 1$  is “heads”,  $X = 2$  is “tails”
- Let  $\mu$  be the probability of heads. If we know  $\mu$ , then:  $P(x = 1|\mu) = \mu$  and  $P(x = 0|\mu) = 1 - \mu$

$$(19) \quad P(x|\mu) = \mu^x(1-\mu)^{1-x} = \begin{cases} \mu & \text{if } x = 1 \\ 1 - \mu & \text{if } x = 0 \end{cases}$$

- This is called the *Bernoulli* distribution

$$(20) \quad E[x] = \mu$$

$$(21) \quad E[(x-\mu)^2] = \mu(1-\mu)$$

- So, suppose we conducted many Bernoulli trials (e.g., flip a coin) and we want to estimate  $\mu$
- **Method: Maximum Likelihood**

$$(22) \quad p(\mathcal{D}|\mu) = \prod_{n=1}^N p(x_n|\mu)$$

$$(23) \quad = \prod_{n=1}^N \mu^{x_n}(1-\mu)^{1-x_n}$$

– Maximize : (*What trick should we use?*)

$$(24) \quad \mathcal{L} = \sum_{n=1}^N x_n \ln \mu + (1 - x_n) \ln(1 - \mu)$$

$$(25) \quad \frac{\partial \mathcal{L}}{\partial \mu} = 0 = \frac{1}{\mu} \sum_{n=1}^N x_n - \frac{1}{1 - \mu} \sum_{n=1}^N (1 - x_n)$$

$$(26) \quad 0 = \frac{(1 - \mu) \sum_{n=1}^N x_n - \mu \sum_{n=1}^N (1 - x_n)}{\mu(1 - \mu)}$$

$$(27) \quad 0 = \sum_{n=1}^N x_n - \mu \sum_{n=1}^N x_n - \mu \sum_{n=1}^N 1 + \mu \sum_{n=1}^N x_n$$

$$(28) \quad 0 = \sum_{n=1}^N x_n - \mu N$$

$$(29) \quad \mu = \frac{1}{N} \sum_{n=1}^N x_n = \frac{m}{N}$$

where  $m$  is the number of successful trials.

- So, if we flip a coin 1 time and get heads, then  $\mu = 1$  and probability of getting tails is 0. *Would you believe that? We need a prior!*
- Look at several independent trials. Consider  $N = 3$  and  $m = 2$  ( $N$  is number of trials,  $m$  is number of successes) and look at all ways to get 2 H and 1 T:
  - \* H H T  $\rightarrow \mu\mu(1 - \mu) = \mu^2(1 - \mu)$
  - \* H T H  $\rightarrow \mu(1 - \mu)\mu = \mu^2(1 - \mu)$
  - \* T H H  $\rightarrow (1 - \mu)\mu\mu = \mu^2(1 - \mu)$
- $\binom{3}{2} \mu^2(1 - \mu) \rightarrow \binom{N}{m} \mu^m(1 - \mu)^{N-m} = \frac{N!}{(N-m)!m!} \mu^m(1 - \mu)^{N-m}$
- This is the Binomial Distribution, gives the probability of  $m$  observations of  $x = 1$  out of  $N$  independent trials
- So, what we saw is that we need a prior. We want to incorporate our prior belief. Let us place a prior on  $\mu$

$$(30) \quad \text{Beta}(\mu|a, b) = \frac{\Gamma(a+b)}{\Gamma(a)\Gamma(b)} \mu^{a-1} (1 - \mu)^{b-1}$$

$$(31) \quad E[\mu] = \frac{a}{a+b}$$

$$(32) \quad \text{Var}[\mu] = \frac{ab}{(a+b)^2(a+b+1)}$$

– Note:  $\Gamma(x) = \int_0^\infty u^{x-1} e^{-u} du$  and when  $x$  is an integer, then it simplifies to  $x!$

- Calculation of the posterior, Take  $N = m + l$  observations:

$$(33) \quad p(\mu|m, l, a, b) \propto \text{Bin}(m, l|\mu) \text{Beta}(\mu|a, b)$$

$$(34) \quad \propto \mu^m (1 - \mu)^l \mu^{a-1} (1 - \mu)^{b-1}$$

$$(35) \quad = \mu^{m+a-1} (1 - \mu)^{l+b-1}$$

- What does this look like? Beta:  $a \leftarrow m + a$ ,  $b \leftarrow l + b$
- So, what's the posterior?

$$(36) \quad p(\mu|m, l, a, b) = \frac{\Gamma(m + a + l + b)}{\Gamma(m + a) \Gamma(l + b)} \mu^{m+a-1} (1 - \mu)^{l+b-1}$$

- **Conjugate Prior Relationship:** When the posterior is the same form as the prior
- Now we can maximize the (log of the) posterior:

$$(37) \quad \max_{\mu} (m + a - 1) \ln \mu + (l + b - 1) \ln(1 - \mu)$$

$$(38) \quad \frac{\partial \mathcal{L}}{\partial \mu} = 0 = \frac{m + a - 1}{\mu} - \frac{l + b - 1}{1 - \mu}$$

$$(39) \quad = (1 - \mu)(m + a - 1) - \mu(l + b - 1)$$

$$(40) \quad = (m + a - 1) - \mu(m + a - 1) - \mu(l + b - 1)$$

$$(41) \quad \mu = \frac{m + a - 1}{m + a + l + b - 2}$$

- This is the MAP solution. *So, what happens now when you flip one heads, two heads, etc.?*
- Discuss online updating of the prior. Eventually the data takes over the prior.