

Lecture 3 Jan 28th

Estimation //

Process for finding unknown parameters that maximize likelihood of underlying distribution.

Common Assumptions:

- Each sample is independent
- Binomial θ is uniform(0,1).

Binomial //

$$p(c_i, \theta) = p(c_i | \theta)p(\theta) = p(\theta | c_i)p(c_i)$$

Thus

$$p(\theta | c_i) = \frac{p(c_i | \theta)p(\theta)}{p(c_i)}$$

So suppose $n=1000$ with 750 heads.

$$p(c_i | \theta)p(\theta) = \prod p(c_i | \theta) = \theta^{750}(1-\theta)^{250}$$

And

$$p(c_i) = \int_0^1 p(c_i | \theta)d\theta = \int_0^1 \theta^{750}(1-\theta)^{250} d\theta = 2$$

$$p(\theta | c_i) = \frac{1}{2} \theta^{750}(1-\theta)^{250}$$

$$\text{Argmax}_{\theta} p(\theta | c_i) = \frac{H}{N} - \text{heads}$$

In General //

Use Baye's Rule:

$$p(\text{model} | \text{data}) = \frac{p(\text{data} | \text{model})p(\text{model})}{p(\text{data})}$$

likelihood

Prior

posterior

evidence

Posterior Dist //

Knowledge of the model based on both the data and the prior.

Likelihood Dist //

Likelihood of sample data assuming model is correct.

Prior Dist //

Assumptions about the model prior to observing the data.

Evidence //

Used in model selection. Could normalize posterior PDF.

2 Ways of Choosing Estimators

Maximum A posteriori

Let $\hat{\theta}$ be the optimal parameter; choose

$$\hat{\theta} = \arg \max_{\theta} p(\theta | D)$$

$$= \arg \max_{\theta} p(D|\theta) p(\theta)$$

Notice $p(D)$, the denominator in Baye's Rule isn't included since it's often just a constant.

Maximum Likelihood

When no prior information is known about θ , we can ignore $p(\theta)$ term and choose,

$$\hat{\theta} = \arg \max_{\theta} p(D|\theta)$$

Tip: Use the negative log function on MAP or ML to easily find the derivative of likelihood functions.

Gaussian Distribution (Normal)

These samples are distributed with $\bar{M} \in \mathbb{R}^n$, $\Sigma \in \mathbb{R}^{n \times n}$.

$$\sim N(\mu, \Sigma)$$

Thus the likelihood func. is,

$$p(D | \mu, \Sigma) = \frac{1}{\sqrt{(2\pi)^n |\Sigma|}} \cdot \exp\left(-\frac{1}{2} (\mathbf{x}_i - \mu)^T \Sigma^{-1} (\mathbf{x}_i - \mu)\right)$$

Then the negative log-likelihood func is,

$$L(D | \mu, \Sigma) = \frac{n}{2} \ln 2\pi + \frac{n}{2} \ln |\Sigma| + \frac{1}{2} \sum (\mathbf{x}_i - \mu)^T \Sigma^{-1} (\mathbf{x}_i - \mu)$$

By solving $\frac{\partial L}{\partial \mu} = 0$ and $\frac{\partial L}{\partial \Sigma}$, we get

$$\mu^* = \frac{1}{n} \sum \mathbf{x}_i$$

$$\Sigma^* = \frac{1}{n} \sum (\mathbf{x}_i - \mu^*)^T (\mathbf{x}_i - \mu^*)^T$$

MAP Non-linear Regression

Now we use a more accurate model for non-linear regression; we account for noise in the data.

$$y = \bar{W}^T b(x) + n$$

where \bar{W} is vector of weights, $b(x)$ is vector with all basis functions and

$$n \sim N(0, \sigma^2)$$

Likelihood of Gaussian

Thus $E(Y) = \vec{w}^T b(x)$

$$\text{Var}(Y) = \sigma^2$$

$$p(\vec{y} | \vec{w}, \vec{x}) = \prod_i \frac{1}{\sqrt{2\pi\sigma^2}} \exp\left(-\frac{1}{2} \frac{(y_i - \vec{w}^T b(x))^2}{\sigma^2}\right)$$
$$= \left(\frac{1}{\sqrt{2\pi\sigma^2}}\right)^n \exp\left(-\frac{1}{2\sigma^2} \sum (y_i - \vec{w}^T b(x))^2\right)$$

Let $\vec{w} \sim N(0, \alpha I)$, (assumption of normal)

$$p(\vec{w} | \vec{x}) = p(\vec{w}) = \frac{1}{\Gamma} \frac{1}{\sqrt{(2\pi\alpha)^d}} e^{-\frac{\|\vec{w}\|^2}{2\alpha}}$$
$$= \frac{1}{(2\pi\alpha)^{\frac{d}{2}}} e^{-\frac{1}{2\alpha} \vec{w}^T \vec{w}}$$

Thus to optimize for \vec{w}^* .

$$p(\vec{w} | \vec{y}, \vec{x}) = \frac{p(\vec{y} | \vec{w}, \vec{x}) p(\vec{w} | \vec{x})}{p(\vec{y} | \vec{x})} \quad \text{by Bayes Rule}$$

Finally

$$L(\vec{w}) = \frac{1}{2} \ln 2\pi\sigma^2 + \frac{1}{2\sigma^2} \sum (y_i - \vec{w}^T b(x))^2$$
$$+ \frac{d}{2} \ln (2\pi\alpha) + \frac{1}{2\alpha} \|\vec{w}\|^2$$
$$+ \ln p(\vec{y} | \vec{x})$$
$$\approx \frac{1}{2\sigma^2} \sum (y_i - \vec{w}^T b(x))^2 + \frac{1}{2\alpha} \|\vec{w}\|^2$$

Since we don't care about the rest when we optimize / derive.

Note: above equation is exactly the same as the least squares method with regularization.

