

TAM 598 Lecture 14 :

Bayesian Linear Regression

Announcements:

- HW 4 covers lectures 13-16; due on Mar 31 (Mon)
- No class Monday Mar 24th

Today:

(1) Max Likelihood Estimation (MLE)

(2) Max A Posteriori Estimation (MPE)

(3) Bayesian Linear Regression

Least Squares Linear Regression - "traditional"

Can we interpret it from a more "modern" Bayesian or probabilistic way?

observations

$$\underline{x}_{1:n}$$

outputs

$$\underline{y}_{1:n}$$

generalized linear model w/ m basis functions:

$$y(\underline{x}; \underline{w}) = \sum_{j=1}^m w_j \underbrace{\phi_j(\underline{x})}_{\substack{\text{basis} \\ \text{weights}}} = \underline{w}^T \underline{\phi}(\underline{x})$$

design matrix
 $\phi_{jj} = \phi_j(\underline{x}_j)$

model the measurement process using a likelihood function

$$\underline{y}_{1:n} \mid (\underline{x}_{1:n}, \underline{w}) \sim p(\underline{y}_{1:n} \mid (\underline{x}_{1:n}, \underline{w}))$$

assume that outcomes of single measurements are Gaussian, with mean $\underline{w}^T \underline{\phi}(\underline{x})$ and noise variance σ^2

$$p(y_i \mid \underline{x}_i, \underline{w}, \sigma^2) = N(y_i \mid y(\underline{x}_i, \underline{w}), \sigma^2) = \underline{N}(y_i \mid \underline{w}^T \underline{\phi}(\underline{x}_i), \sigma^2)$$

notation: $N(y|\mu, \sigma^2)$ is a pdf

$$p(y) = \underbrace{(2\pi\sigma^2)^{-1/2}}_{\text{constant}} \exp \left\{ -\frac{(y-\mu)^2}{2\sigma^2} \right\}$$

for independent measurements, the likelihood of the data factorizes

$$\underbrace{P(\underline{y}_{1:n} | \underline{x}_{1:n}, \underline{w})}_{\text{likelihood}} = \prod_{i=1}^n P(y_i | x_i, w)$$

$$= \prod_{i=1}^n N(y_i | \underline{w}^T \underline{\Phi}(x_i), \sigma^2)$$

⋮

$$\cancel{*} = (2\pi\sigma^2)^{-n/2} \exp \left\{ -\frac{1}{2\sigma^2} \| \underline{y}_{1:n} - \underline{\Phi} \underline{w} \| \right\}$$

^{n × m}
design matrix

Today:

(1) Max Likelihood Estimation (MLE)

- assume measurement outcomes are gaussian distributed

- find weights ω and variance σ^2 that maximizes likelihood of measuring the observed data

(2) Max A Posteriori Estimation (MPE)

(3) Bayesian Linear Regression

I How do we find the parameters? Maximum likelihood for weights \underline{w} , and σ^2

$$\max_{\underline{w}, \sigma^2} \log P(\underline{y}_{1:n} | \underline{x}_{1:n}, \underline{w})$$

$$\max_{\underline{w}, \sigma^2} \left\{ -\frac{n}{2} \log(2\pi) - \frac{n}{2} \log \sigma^2 - \frac{1}{2\sigma^2} \left\| \underline{y}_{1:n} - \underline{\Phi} \underline{w} \right\|^2 \right\}$$

neg of sum of square errors

So: maximizing likelihood wrt \underline{w} is equivalent to minimizing sum of squares! And our weights should therefore satisfy

$$\underline{\Phi}^T \underline{\Phi} \underline{w} = \underline{\Phi}^T \underline{y}_{1:n}$$

AND also we can estimate σ^2 as well via max likelihood

$$0 = \frac{d}{d\sigma^2} \{ \} \Rightarrow$$

$$\sigma^2 = \frac{\|\underline{\Phi} \underline{w} - \underline{y}_{1:n}\|^2}{n}$$

Now: we can incorporate uncertainty σ^2 when making predictions

point predictive distribution : $P(y | \underline{x}, \underline{w}, \sigma^2) = N(y | \underline{w}^\top \phi(\underline{x}), \sigma^2)$

the measured output y is normal distributed around the model (least squares) prediction, with variance σ^2

Today:

(1) Max Likelihood Estimation (MLE)

(2) Max A Posteriori Estimation (MPE)

- ✓ - assume measurement outcomes are gaussian distributed as before
- ✓ - assume gaussian prior on weights $p(\underline{w})$ with zero mean, variance α^2
- ✓ - assume measurement uncertainty σ^2
- ✓ - find the weights \underline{w} that maximize the prob. of the Bayesian posterior
- ✓ - helps to avoid overfitting

(3) Bayesian Linear Regression

II How do we find the parameters? Maximum a posteriori estimates

here: maximize the log prob of the posterior rather than the likelihood. Helps avoid overfitting

$$\text{again: } y(\underline{x}; \underline{w}) = \underline{w}^T \underline{\phi}(\underline{x})$$

$$\underline{y}_{1:n} \mid \underline{x}_{1:n}, \underline{w}, \sigma^2 \sim N(\underline{w}^\top \underline{\phi}(\underline{x}), \sigma^2)$$

 now assume

X new: uncertainty in model parameters described by a prior $w \sim p(w)$

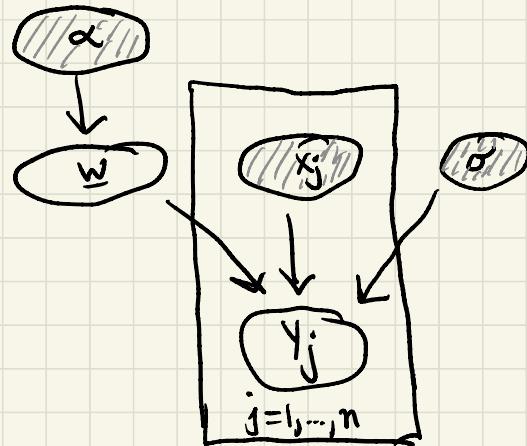
eg) gaussian prior on weights $p(\underline{w} | \alpha) = N(\underline{w} | \underline{\theta}, \alpha^{-1} \mathbb{I})$

mean values

covariance matrix

$$= \left(\frac{\alpha}{2\pi} \right)^{\frac{m}{2}} \exp \left\{ -\frac{\alpha}{2} \|w\|^2 \right\}$$

graphically



posterior from Bayes' Rule

$$p(\underline{w} | \underline{x}_{1:n}, \underline{y}_{1:n}, \sigma, \alpha) =$$

$$\frac{P(\underline{y}_{1:n} | \underline{x}_{1:n}, \underline{w}, \sigma, \alpha) P(\underline{w} | \alpha)}{\int P(\underline{y}_{1:n} | \underline{x}_{1:n}, \underline{w}', \sigma, \alpha) P(\underline{w}' | \alpha) d\underline{w}'}$$

our state of knowledge about \underline{w} , after we see the data

Now assuming σ and α are known.

a point estimate of \underline{w} is

$$\underline{w}_{MPE} = \operatorname{argmax}_{\underline{w}} \underbrace{P(\underline{y}_{1:n} | \underline{x}_{1:n}, \underline{w}, \sigma^2)}_{\text{likelihood}} \underbrace{p(\underline{w} | \alpha)}_{\text{prior}}$$

For gaussian likelihood and weight prior:

$$\log p(\underline{w} | \underline{x}_{1:n}, \underline{y}_{1:n}, \sigma^2, \alpha) = \left\{ -\frac{1}{2\sigma^2} \|\underline{y}_{1:n} - \underline{w}\|_2^2 - \frac{\alpha}{2} \|\underline{w}\|_2^2 \right\}$$

maximizing:

$$\nabla_{\underline{w}} \{ \cdot \} = 0 \Rightarrow$$

$$\boxed{\underline{w}_{MPE} = \frac{1}{\sigma^2} \left[\frac{1}{\sigma^2} \underline{\Phi}^T \underline{\Phi} + \alpha \underline{I} \right]^{-1} \underline{\Phi}^T \underline{y}_{1:n}}$$

Today:

(1) Max Likelihood Estimation (MLE)

(2) Max A Posteriori Estimation (MPE)

(3) Bayesian Linear Regression

- just like (2) above, but now work with the full posterior distribution, not just a point estimate
- allows us to separate epistemic and aleatoric uncertainty

OR another way to see this: our posterior

$$p(\underline{w} | \underline{x}_{1:n}, \underline{y}_{1:n}, \underline{\sigma^2}, \alpha) \propto \exp \left\{ -\frac{1}{2\sigma^2} \| \underline{\Phi} \underline{w} - \underline{y}_{1:n} \|^2 - \frac{\alpha}{2} \| \underline{w} \|^2 \right\}$$

algebra: Rewrite in general form of a gaussian

$$\propto \exp \left\{ -\frac{1}{2} (\underline{w} - \underline{m})^\top \underline{S}^{-1} (\underline{w} - \underline{m}) \right\}$$

where $\underline{m} = \frac{1}{\sigma^2} \underline{\Phi}^T \underline{y}_{1:n}$ is the mean vector

- If $\alpha \rightarrow 0$:
OLS solution
- If $\sigma^2 \rightarrow 0$, data contribution dominates prior, more weight in \underline{w}
 $\underline{S} = \left[-\frac{1}{\sigma^2} \underline{\Phi}^T \underline{\Phi} + \alpha \underline{I} \right]^{-1}$ is the covariance matrix
 - from likelihood (data contrib)
 - from prior, regularizes est of \underline{w}

III

How do we find the parameters? Bayesian Linear Regression

similar to II :

again: $\underline{y}(\underline{x}; \underline{w}) = \underline{w}^T \underline{\phi}(\underline{x})$

$$\underline{Y}_{1:n} \mid \underline{x}_{1:n}, \underline{w}, \sigma^2 \sim N(\underline{w}^T \underline{\phi}(\underline{x}), \sigma^2)$$

\underline{w} , σ^2

~~IV~~ uncertainty in model parameters described by a prior $\underline{w} \sim p(\underline{w})$

but no point estimate. Keep the posterior distribution and work with it directly. Why? Can now quantify the epistemic uncertainty arising from limited # of observations.

posterior is gaussian : $p(\underline{w} \mid \underline{x}_{1:n}, \underline{y}_{1:n}, \sigma^2, \alpha) = N(\underline{w} \mid \underline{m}, \underline{S})$

where $\underline{S} = \left(\frac{1}{\sigma^2} \underline{\Phi}^T \underline{\Phi} + \alpha \underline{I} \right)^{-1}$; $\underline{m} = \frac{1}{\sigma^2} \underline{S} \underline{\Phi}^T \underline{y}_{1:n}$

posterior predictive distribution: what can we say about y at some new \underline{x} after seeing the data?

$$P(\underline{y} | \underline{\underline{x}}, \underline{\underline{x}}_{1:n}, \underline{\underline{y}}_{1:n}, \sigma^2, \alpha) = \int P(\underline{y} | \underline{\underline{x}}, \underline{\underline{w}}, \sigma^2) P(\underline{\underline{w}} | \underline{\underline{x}}_{1:n}, \underline{\underline{y}}_{1:n}, \sigma^2, \alpha) d\underline{\underline{w}}$$

for all gaussian priors, this is analytically available

$$P(\underline{\underline{y}} | \underline{\underline{x}}, \underline{\underline{x}}_{1:n}, \underline{\underline{y}}_{1:n}, \sigma^2, \alpha) = N(\underline{\underline{y}} | m(\underline{\underline{x}}), s^2(\underline{\underline{x}}))$$

where $m(\underline{\underline{x}}) = \underline{\underline{m}}^\top \underline{\phi}(\underline{\underline{x}})$

$$s^2(\underline{\underline{x}}) = \underline{\underline{\phi}}^\top(\underline{\underline{x}}) \underline{\underline{\Sigma}} \underline{\underline{\phi}}(\underline{\underline{x}}) + \sigma^2$$

epistemic uncertainty measurement noise