Introduction to Machine Learning - CS725
Instructor: Prof. Ganesh Ramakrishnan
Lecture 5 - Linear Regression - Bayesian
Inference and Regularization

# Building on questions on Least Squares Linear Regression

1. Is there a probabilistic interpretation?
   - Gaussian Error, Maximum Likelihood Estimate
2. Addressing overfitting
   - Bayesian and Maximum Aposteriori Estimates, Regularization
3. How to minimize the resultant and more complex error functions?
   - Level Curves and Surfaces, Gradient Vector, Directional Derivative, Gradient Descent Algorithm, Convexity, Necessary and Sufficient Conditions for Optimality

# Recap: Illustration through a Simple Coin Tossing Example: Maximum Likelihood Estimation vs. Bayesian Estimation

Suresh likes to toss coins. One day he decided to count the number of heads and tails in his coin tosses. Here is what he found. After tossing 1000 times (it took him a hours, but he likes to toss coins), he found that the coin landed on heads 400 times and tails 600 times. His reflection: If I were to toss the coin once more time, what is the probability that I get a heads?

# Recap: MLE estimate for Coin Tossing

- We restate Suresh's problem as the MLE of the probability of getting a head. This is the value of $p$ which maximizes the likelihood of observing 400 heads as outcomes.

$$\hat{p} = \underset{p}{\text{argmax}} \; {}^{1000}C_{400}p^{400}(1-p)^{600}$$

- $\hat{p} = 0.4$ as we had intuitively guessed. In general, the value of $p$ which maximises the likelihood of observing $h$ heads, given $n$ coin tosses is.

$$\hat{p} = \underset{p}{\text{argmax}} \; {}^{n}C_{h}p^{h}(1-p)^{n-h}$$

# Bayesian Inference/Estimation

Suresh now brings a newly minted coin to toss. He *believes* that the coin is fair and heads and tails are equally likely outcomes (since the coin is not worn out). Now like always he flips the coin 4 times, and finds out that heads appeared all the 4 times.

1. Is the MLE estimate $\hat{p} = 1$ intuitive? Is tails *improbable*?
2. Is there a way that Suresh could update his *belief* about the coin.

# Bayesian Inference

- $H$: One of few competing hypotheses whose probability may be affected by observed data.
- $\Pr(H)$: The (prior) probability of $H$ before data $\mathcal{D}$ is observed. This indicates one's previous *belief* in the hypothesis.
- The evidence $\mathcal{D}$: New data that were not used in computing the prior probability

$$p(H \mid \mathcal{D}) \propto p(\mathcal{D} \mid H)\, p(H)$$

# Conjugate Prior

Let $\mathcal{D} \mid H$ follow a distribution $d_1$ and $H$ follow a distribution $d_2$.
The distribution $d_2$ is the conjugate prior of $d_1$ if the distribution of $\Pr(H \mid \mathcal{D})$ follows the distribution $d_2$.
Some Examples:

1. Bernoulli & Binomial ($d_1$) - Beta ($d_2$)
2. Geometric - Beta
3. Categorical - Dirichlet
4. Multinomial - Dirichlet
5. Poisson - Gamma
6. Normal - Inverse Gamma

Recall form of binomial: $^nC_h \, p^h \, (1-p)^{n-h}$ $(n-h=t)$

$^{h+t}C_h \, p^h \, (1-p)^t$

Let $\mathcal{D} \mid H$ follow a distribution $Ber(p)$ ($p$ is probability of heads) and $p$ follow a distribution $Beta(p; \alpha, \beta) \sim \dfrac{p^{(\alpha-1)}(1-p)^{(\beta-1)}}{B(\alpha, \beta)}$,

- *The beta normalization function:*

$P(H) = B(\alpha, \beta) = \displaystyle\int_{p=0}^{1} p^{(\alpha-1)}(1-p)^{(\beta-1)} dp = \dfrac{\Gamma(\alpha)\Gamma(\beta)}{\Gamma(\alpha+\beta)}$, where $\Gamma(.)$

behaves like the factorial function: $\Gamma(n) = (n-1)!$ if $n \in \mathbb{Z}^+$

If $\alpha, \beta \in \mathbb{Z}^+$ $B(\alpha, \beta) = \dfrac{1}{^{\alpha+\beta-1}C_{\alpha-1} \times \beta}$

is $\dfrac{1}{b(\alpha, \beta)} = \;^{\alpha+\beta-1}C_{\alpha-1} \times \beta$ $\longrightarrow$ Similar to $^{h+t}C_h$

# The Beta Conjugate Prior for Bernoulli/Binomial

$p$ is a param

Let $\mathcal{D} \mid H$ follow a distribution $Ber(p)$ ($p$ is probability of heads) and $p$ follow a distribution $Beta(p; \alpha, \beta) \sim \frac{p^{(\alpha-1)}(1-p)^{(\beta-1)}}{B(\alpha,\beta)}$,

- $\Pr(H \mid \mathcal{D}) = \Pr(p \mid \mathcal{D}) = \dfrac{p(x_1 \ldots x_n \mid p)\, p(p)}{\int p(x_1 \ldots x_n \mid p)\, p(p)}$

we know from conjugacy that this MUST be a Beta distribution

$$= \frac{{}^nC_h\, p^h (1-p)^{n-h} \ast p^{(\alpha-1)}(1-p)^{\beta-1}}{\int_p \cdots \cdots}$$

$$= \frac{p}{}\; p^{h+\alpha-1}(1-p)^{n-h+\beta-1} \ast \cdots$$

$$= Beta\left(p;\; \alpha+h,\; \beta+n-h\right)$$

# The Beta Conjugate Prior for Bernoulli/Binomial

Let $\mathcal{D} \mid H$ follow a distribution $Ber(p)$ ($p$ is probability of heads) and $p$ follow a distribution $Beta(p; \alpha, \beta) \sim \frac{p^{(\alpha-1)}(1-p)^{(\beta-1)}}{B(\alpha, \beta)}$,

- $\Pr(H \mid \mathcal{D}) = \Pr(p \mid \mathcal{D}) = \dfrac{\Pr(\mathcal{D} \mid p)\Pr(p)}{\displaystyle\int_q \Pr(\mathcal{D} \mid q)\Pr(q)}$

$$= \frac{{}^nC_h p^h(1-p)^{n-h}\frac{\Gamma(\alpha+\beta)}{\Gamma(\alpha)\Gamma(\beta)}p^{(\alpha-1)}(1-p)^{(\beta-1)}}{\displaystyle\int_q {}^nC_h q^h(1-q)^{n-h}\frac{\Gamma(\alpha+\beta)}{\Gamma(\alpha)\Gamma(\beta)}q^{(\alpha-1)}(1-q)^{(\beta-1)}}$$

$$\propto p^{\alpha+h-1}(1-p)^{\beta+n-h-1} \sim Beta(p; \alpha+h, \beta+n-h)$$

*Added # heads to $\alpha$ &*
*# tails to $\beta$*

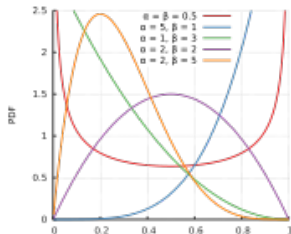# More on the $Beta(\alpha, \beta)$ distribution

1. $\mathbf{E}_{Beta(\alpha,\beta)}[p] = \frac{\alpha}{\alpha+\beta}$ and $\underset{p}{\operatorname{argmax}}\ Beta(p; \alpha, \beta) = \frac{\alpha - 1}{\alpha + \beta - 2}$
   (the mode of the distribution)

   *[handwritten annotations: "Mean" pointing to $\mathbf{E}$; circle around $\frac{\alpha}{\alpha+\beta}$; "model" pointing to argmax; $\approx \frac{\#\ prior\ heads}{\#\ prior\ tosses}$]*

2. $Beta(1, 1)$ is the uniform distribution!

3. Is the conjugate prior pdf for the Bernoulli, binomial, negative binomial and geometric distributions and has the following pdf:

# The MAP Estimate for Bernoulli/Binoimal

Let $\mathcal{D} \mid H$ follow a distribution $Ber(p)$ ($p$ is probability of heads)
and $p$ follow a distribution $Beta(p; \alpha, \beta) \sim \frac{p^{(\alpha-1)}(1-p)^{(\beta-1)}}{B(\alpha, \beta)}$,

1. *The Maximum Likelihood Estimate:*
   $$\hat{p} = \underset{p}{\operatorname{argmax}} \ {}^nC_h p^h (1-p)^{n-h} = \frac{h}{n}$$

2. *The Maximum a-Posterior (MAP) Estimate:* The mode of the posterior distribution
   $$\tilde{p} = \underset{H}{\operatorname{argmax}} \ \Pr(H \mid \mathcal{D}) = \underset{p}{\operatorname{argmax}} \ \Pr(p \mid \mathcal{D}) = \frac{\alpha + h - 1}{\alpha + \beta + n - 2}$$

Recall: $\underset{p}{\operatorname{argmax}} \ Beta(p; \alpha, \beta) = \frac{\alpha - 1}{\alpha + \beta - 2}$

# The MAP Estimate for Bernoulli/Binoimal

Let $\mathcal{D} \mid H$ follow a distribution $Ber(p)$ ($p$ is probability of heads) and $p$ follow a distribution $Beta(p; \alpha, \beta) \sim \frac{p^{(\alpha-1)}(1-p)^{(\beta-1)}}{B(\alpha,\beta)}$,

1. *The Maximum Likelihood Estimate:*
   $$\hat{p} = \underset{p}{\mathrm{argmax}} \; {}^nC_h p^h (1-p)^{n-h} = \frac{h}{n}$$

2. *The Maximum a-Posterior (MAP) Estimate:* The mode of the posterior distribution
   $$\tilde{p} = \underset{H}{\mathrm{argmax}} \; \Pr(H \mid \mathcal{D}) = \underset{p}{\mathrm{argmax}} \; \Pr(p \mid \mathcal{D})$$
   $$= \underset{p}{\mathrm{argmax}} \; Beta(p; \alpha + h, \beta + n - h) = \frac{\alpha + h - 1}{\alpha + \beta + n - 2}$$

# Case Study Continued

Coming back to the Suresh's case study, he observed 4 heads on 4 tosses, his MLE is

$$\hat{p} = \underset{p}{\operatorname{argmax}} \ ^{4}C_{4}p^{4}(1-p)^{0} = 1$$

If his prior on $p$ was $Beta(p; 3, 3)$, then his posterior will be $Beta(p; 3+4, 3+0) = Beta(p; 7, 3)$ and his MAP estimate will be

$$\hat{p} = \underset{p}{\operatorname{argmax}} \ Beta(p; 7, 3) = \frac{7-1}{7+3-2} = 0.75$$

$\alpha = 3 \quad \beta = 3$

Often, $\alpha, \beta$ are tuned to maximize performance on validation / held-out data

# Summary: Bayesian Inference for Bernoulli

Let $\mathcal{D} \mid H$ follow a distribution $Ber(p)$ ($p$ is probability of heads) and $p$ follow a distribution $Beta(p; \alpha, \beta) \sim \frac{p^{(\alpha-1)}(1-p)^{(\beta-1)}}{B(\alpha,\beta)}$,

1. *The Maximum Likelihood Estimate:*
   $$\hat{p} = \underset{p}{\mathrm{argmax}} \ ^{n}C_{h} p^{h}(1-p)^{n-h} = \frac{h}{n}$$

2. *The Posterior Distribution:*
   $$\Pr(p \mid \mathcal{D}) = Beta(p; \alpha + h, \beta + n - h)$$

3. *The Maximum a-Posterior (MAP) Estimate:* The mode of the posterior distribution
   $$\tilde{p} = \underset{H}{\mathrm{argmax}} \Pr(H \mid \mathcal{D}) = \underset{p}{\mathrm{argmax}} \Pr(p \mid \mathcal{D})$$
   $$= \mathrm{argmax}\, Beta(p; \alpha + h, \beta + n - h) = \frac{\alpha + h - 1}{\alpha + \beta + n - 2}$$

# Illustration of Bayesian Estimation on a continuous valued random variable

# Conjugate Prior for (univariate) Gaussian

- Let $\Pr(X) \sim \mathcal{N}(\mu, \sigma^2)$ and let the data $\mathcal{D} = x_1 \ldots x_m$

- $\mu_{MLE} = \frac{1}{m} \sum_{i=1}^{m} x_i$ and $\sigma^2_{MLE} = \frac{1}{m} \sum_{i=1}^{m} (x_i - \mu_{MLE})^2$

$$\hat{\mu}, \hat{\sigma} = \arg\max_{\mu, \sigma} \prod_i N(x_i | \mu, \sigma)$$

- Suppose you are told that the conjugate prior for the (univariate) normally distributed random variable $X$ in the case that $\sigma^2$ is not a random variable is *(prior belief on $\mu$, NOT on $\sigma^2$)* $\Pr(\mu) = \mathcal{N}(\mu_0, \sigma_0^2)$. Then the **posterior** is?

$$\Pr(H | D) = \Pr(\mu | D) = \mathcal{N}(\mu_n, \sigma_n^2)$$

How will you find $\mu_n$ & $\sigma_n^2$? Trick??

$$p(x_1 \ldots x_n | \mu, \sigma^2) = \left(\frac{1}{\sigma\sqrt{2\pi}}\right)^n \exp\left(\sum_{i=1}^{n} \frac{(x_1 - \mu)^2}{2\sigma^2}\right)$$

$$\mu_{MLE}, \sigma^2_{MLE} = \underset{\mu, \sigma^2}{\text{argmax}} \left(\frac{1}{\sigma\sqrt{2\pi}}\right)^n \exp\left(-\sum_{i=1}^{n} \frac{(x_i - \mu)^2}{2\sigma^2}\right)$$

$$\mu_{MLE} = \underset{\mu}{\text{argmax}} -\sum_{i=1}^{n} (x_1 - \mu)^2 \qquad \left(\text{Intuitively, } \mu_{MLE} \text{ is independent of } \sigma^2\right)$$

Recall, mean mean minimizes sum of squares of deviation $\Rightarrow \mu_{MLE} = \frac{1}{n}\sum_i x_i$

# Conjugate Prior for (univariate) Gaussian

- Let $\Pr(X) \sim \mathcal{N}(\mu, \sigma^2)$ and let the data $\mathcal{D} = x_1 \ldots x_m$
- $\mu_{MLE} = \frac{1}{m} \sum_{i=1}^{m} x_i$ and $\sigma^2_{MLE} = \frac{1}{m} \sum_{i=1}^{m} (x_i - \mu_{MLE})^2$
- Suppose you are told that the conjugate prior for the (univariate) normally distributed random variable $X$ in the case that $\sigma^2$ is not a random variable is $\Pr(\mu) = \mathcal{N}(\mu_0, \sigma_0^2)$. Then the **posterior** is?
- Answer: $\Pr(\mu | x_1 \ldots x_m) = \mathcal{N}(\mu_m, \sigma_m^2)$ such that $\mu_m = \ldots\ldots$ and $\frac{1}{\sigma_m^2} = \ldots..$
- Helpful tip: Product of Gaussians is always a Gaussian

$$\text{①} \Pr(\mu) = \frac{1}{\sqrt{2\pi\sigma_0^2}} exp\left(\frac{-(\mu - \mu_0)^2}{2\sigma_0^2}\right) \sim N\left(\mu_0, \sigma_0^2\right)$$

$$\Pr(x_i|\mu; \sigma^2) = \frac{1}{\sqrt{2\pi\sigma^2}} exp\left(\frac{-(x_i - \mu)^2}{2\sigma^2}\right) \sim N\left(\mu, \sigma^2\right)$$

$$\text{②} \Pr(\mathcal{D}|\mu) = \left(\frac{1}{\sqrt{2\pi\sigma^2}}\right)^m exp\left(\frac{-1}{2\sigma^2}\sum_{i=1}^{m}(x_i - \mu)^2\right)$$

$$\Pr(\mu|\mathcal{D}) \propto \Pr(\mathcal{D}|\mu)\Pr(\mu) = \text{①} \wedge \text{②}$$

$$\left(\frac{1}{\sqrt{2\pi\sigma^2}}\right)^m \frac{1}{\sqrt{2\pi\sigma_0^2}} exp\left(\frac{-1}{2\sigma^2}\sum_{i=1}^{m}(x_i - \mu)^2 - \frac{(\mu - \mu_0)^2}{2\sigma_0^2}\right) \propto$$

$$exp\left(\frac{-1}{2\sigma^2}\sum_{i=1}^{m}(x_i - \mu)^2 - \frac{(\mu - \mu_0)^2}{2\sigma_0^2}\right) = exp\left(\frac{-1}{2\sigma_m^2}(\mu - \mu_m)^2\right)$$

$$N\left(\mu_m, \sigma_m^2\right)$$

# Detailed derivation (contd.)

Our reference equality: *(from $P_r(D|M)$)* *(from $P_r(M)$)*

$$exp\left(\frac{-1}{2\sigma^2}\sum_{i=1}^{m}(x_i - \mu)^2 - \frac{(\mu - \mu_0)^2}{2\sigma_0^2}\right) = exp\left(\frac{-1}{2\sigma_m^2}(\mu - \mu_m)^2\right),$$

Matching coefficients of $\mu^2$, we get

$$\frac{-1}{2\sigma^2}\left(m\mu^2\right) - \frac{\mu^2}{2\sigma_0^2} = \frac{-1}{2\sigma_m^2}\mu^2$$

Since we believe that $P_r(M|D)$ follows Normal distr.

Our reference equality:

$$exp\left(\frac{-1}{2\sigma^2}\sum_{i=1}^{m}(x_i - \mu)^2 - \frac{(\mu - \mu_0)^2}{2\sigma_0^2}\right) = exp\left(\frac{-1}{2\sigma_m^2}(\mu - \mu_m)^2\right),$$

Matching coefficients of $\mu^2$, we get

$$\frac{-\mu^2}{2\sigma_m^2} = \frac{-\mu^2}{2}\left(\frac{m}{\sigma^2} + \frac{1}{\sigma_0^2}\right) \Rightarrow$$

Our reference equality:

$$exp\left(\frac{-1}{2\sigma^2}\sum_{i=1}^{m}(x_i - \mu)^2 - \frac{(\mu - \mu_0)^2}{2\sigma_0^2}\right) = exp\left(\frac{-1}{2\sigma_m^2}(\mu - \mu_m)^2\right),$$

Matching coefficients of $\mu^2$, we get

$$\frac{-\mu^2}{2\sigma_m^2} = \frac{-\mu^2}{2}\left(\frac{m}{\sigma^2} + \frac{1}{\sigma_0^2}\right) \Rightarrow \frac{1}{\sigma_m^2} = \frac{1}{\sigma_0^2} + \frac{m}{\sigma^2}$$

Matching coefficients of $\mu$, we get

# Detailed derivation (contd.)

Our reference equality:

$$exp\left(\frac{-1}{2\sigma^2}\sum_{i=1}^{m}(x_i - \mu)^2 - \frac{(\mu - \mu_0)^2}{2\sigma_0^2}\right) = exp\left(\frac{-1}{2\sigma_m^2}(\mu - \mu_m)^2\right),$$

Matching coefficients of $\mu^2$, we get

$$\frac{-\mu^2}{2\sigma_m^2} = \frac{-\mu^2}{2}\left(\frac{m}{\sigma^2} + \frac{1}{\sigma_0^2}\right) \Rightarrow \frac{1}{\sigma_m^2} = \frac{1}{\sigma_0^2} + \frac{m}{\sigma^2}$$

Matching coefficients of $\mu$, we get

$$\frac{2\mu\mu_m}{2\sigma_m^2} = \mu\left(\frac{2\sum_{i=1}^{m}x_i}{2\sigma^2} + \frac{2\mu_0}{2\sigma_0^2}\right) \Rightarrow$$

Our reference equality:

$$exp\left(\frac{-1}{2\sigma^2}\sum_{i=1}^{m}(x_i - \mu)^2 - \frac{(\mu - \mu_0)^2}{2\sigma_0^2}\right) = exp\left(\frac{-1}{2\sigma_m^2}(\mu - \mu_m)^2\right),$$

Matching coefficients of $\mu^2$, we get

$$\frac{-\mu^2}{2\sigma_m^2} = \frac{-\mu^2}{2}\left(\frac{m}{\sigma^2} + \frac{1}{\sigma_0^2}\right) \Rightarrow \frac{1}{\sigma_m^2} = \frac{1}{\sigma_0^2} + \frac{m}{\sigma^2}$$

Matching coefficients of $\mu$, we get

$$\frac{2\mu\mu_m}{2\sigma_m^2} = \mu\left(\frac{2\sum_{i=1}^{m}x_i}{2\sigma^2} + \frac{2\mu_0}{2\sigma_0^2}\right) \Rightarrow \mu_m = \sigma_m^2\left(\frac{\sum_{i=1}^{m}x_i}{\sigma^2} + \frac{\mu_0}{\sigma_0^2}\right) \text{ or }$$

$$\mu_m = \sigma_m^2\left(\frac{m\hat{\mu}_{ML}}{\sigma^2} + \frac{\mu_0}{\sigma_0^2}\right) \Rightarrow$$

Our reference equality:

$$exp\left(\frac{-1}{2\sigma^2}\sum_{i=1}^{m}(x_i - \mu)^2 - \frac{(\mu - \mu_0)^2}{2\sigma_0^2}\right) = exp\left(\frac{-1}{2\sigma_m^2}(\mu - \mu_m)^2\right),$$

Matching coefficients of $\mu^2$, we get

$$\frac{-\mu^2}{2\sigma_m^2} = \frac{-\mu^2}{2}\left(\frac{m}{\sigma^2} + \frac{1}{\sigma_0^2}\right) \Rightarrow \frac{1}{\sigma_m^2} = \frac{1}{\sigma_0^2} + \frac{m}{\sigma^2}$$

Matching coefficients of $\mu$, we get

$$\frac{2\mu\mu_m}{2\sigma_m^2} = \mu\left(\frac{2\sum_{i=1}^{m}x_i}{2\sigma^2} + \frac{2\mu_0}{2\sigma_0^2}\right) \Rightarrow \mu_m = \sigma_m^2\left(\frac{\sum_{i=1}^{m}x_i}{\sigma^2} + \frac{\mu_0}{\sigma_0^2}\right) \text{ or}$$

$$\mu_m = \sigma_m^2\left(\frac{m\hat{\mu}_{ML}}{\sigma^2} + \frac{\mu_0}{\sigma_0^2}\right) \Rightarrow \mu_m = \left(\frac{\sigma^2}{m\sigma_0^2 + \sigma^2}\mu_0\right) + \left(\frac{m\sigma_0^2}{m\sigma_0^2 + \sigma^2}\hat{\mu}_{ML}\right)$$