

Gradient Descent

- Initialisation: $w \leftarrow w_0$
- Repeat until convergence: $\|\nabla_w E\| < \epsilon$

$$w_{t+1} \leftarrow w_t - \eta \nabla_w E \quad \eta: \text{learning rate}$$

$$E(w, D) = \sum_{i=1}^n (w^T x_i - y_i)^2$$

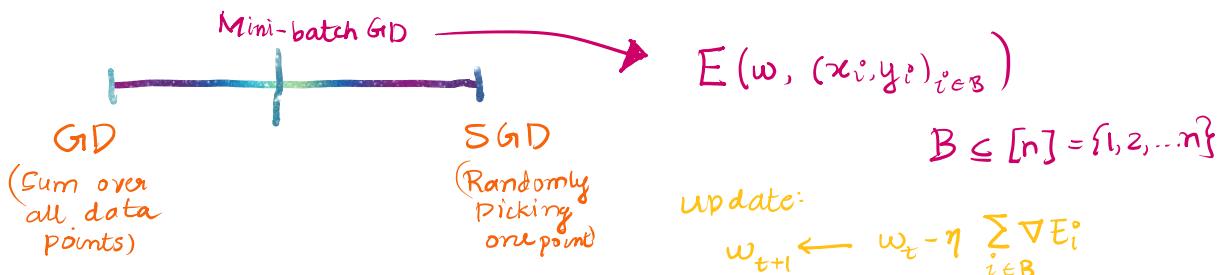
$$\nabla_w E = \sum_{i=1}^n \nabla_w E_i$$

Gradient descent is excellent in terms of accuracy but expensive in computation

Stochastic Gradient Descent

update state: $w_{t+1} \leftarrow w_t - \nabla_w E(w, x_i, y_i)$, (i is randomly chosen) $\rightarrow \nabla_w ((w^T x_i - y_i)^2)$

Faster computation than Gradient descent

MLE (recap):

Let $D = \{(x_i, y_i)\}_{i \in [n]}$ and θ be parameter to be estimated.

We are interested in $\arg \max_{\theta} P(D | \theta) = \theta_{MLE}$
↳ likelihood function

Recall the coin toss example here.

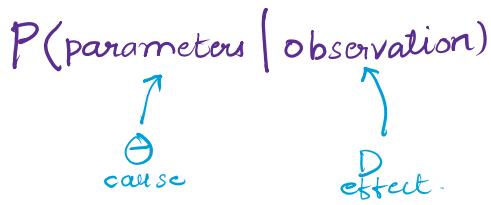
$$\text{In the coin toss, } \theta_{MLE} = \frac{1}{n} \sum_{j=1}^n y_j$$

Now, what if we had a prior knowledge regarding θ ?

Eg: In coin toss, we know a-priori that θ has a normal-like distribution

Can we update / incorporate this idea? This leads to the notion of Maximum A posteriori Estimate.

Maximum A posteriori Estimate (MAP):



Bayesian inference:

Prior ↓

likelihood ↘

$$P(\theta \mid D) = \frac{P(D \mid \theta) P(\theta)}{P(D)}$$

We call $P(\theta \mid D)$ posterior belief

$$\Theta_{\text{MAP}} \in \underset{\theta}{\operatorname{argmax}} P(\theta \mid D) = \underset{\theta}{\operatorname{argmax}} [P(D \mid \theta) P(\theta)]$$

Now, $\log P(\theta \mid D) = \log P(D \mid \theta) + \log P(\theta)$

$$\Rightarrow \Theta_{\text{MAP}} \in \underset{\theta}{\operatorname{argmax}} [\log P(D \mid \theta) + \log P(\theta)]$$

If $P(\theta)$ is independent of θ (say uniform) then $\Theta_{\text{MAP}} = \Theta_{\text{MLE}}$
 i.e., $P(\theta)$ is a constant.

Ex: Likelihood at observing k heads in n tosses.

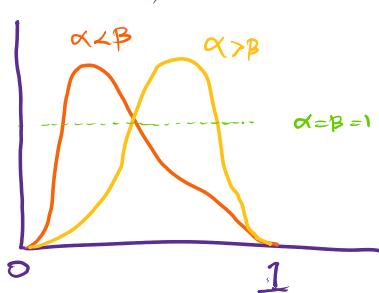
$$P(D \mid \theta) = \binom{n}{k} \theta^k (1-\theta)^{n-k} \quad \text{Bin}(n, \theta)$$

Here, $\Theta_{\text{MLE}} = k/n$

Beta distribution

pmf: $P(\theta) = \frac{1}{C} \theta^{\alpha-1} (1-\theta)^{\beta-1}$

If $\alpha = \beta = 1$, $P(\theta)$ becomes constant, we've a uniform distribution.



- Beta includes a large family of distributions in $[0, 1]$
- Beta is a conjugate prior at binomial distribution.

Conjugate prior (CP):

Let $P(D \mid \theta) \sim d_1$, $P(\theta) = d_2$

Then, $P(\theta)$ is a CP of $P(D \mid \theta)$
 if $P(\theta \mid D) \sim d_2$

$$\begin{aligned}
 P(\theta \mid D) &\propto \underbrace{\theta^k}_{P(D \mid \theta)} \underbrace{(1-\theta)^{n-k}}_{P(\theta)} \\
 &\propto \theta^{k+\alpha-1} (1-\theta)^{n-k+\beta-1} \sim \text{Beta}(k+\alpha, n-k+\beta)
 \end{aligned}$$

$$\underset{\theta}{\operatorname{argmax}} P(\theta \mid D) = \underset{\theta}{\operatorname{argmax}} \log P(\theta \mid D)$$

$$= \text{const.} + (k+\alpha-1) \log \theta + (n-k+\beta-1) \log (1-\theta)$$

$$\underset{\theta}{\sim} \text{Beta}(k+\alpha-1, n-k+\beta-1) \quad (\text{If } \alpha = \beta = 1, \Theta_{\text{MAP}} = \Theta_{\text{MLE}})$$

$$= \text{const.} + (k+\alpha-1) \log \theta + (n-k+p-1) \log(1-\theta)$$

$$\hat{\theta}_{\text{MAP}} = \frac{k+\alpha-1}{n+\alpha+\beta-2}$$

Θ_{MLE}

(If $\alpha=\beta=1$, $\hat{\theta}_{\text{MAP}} = \hat{\theta}_{\text{MLE}}$
uniform \rightarrow no prior info.)

Conjugate prior examples

1. Bernoulli/Binomial \leftrightarrow Beta
2. Geometric \leftrightarrow Beta
3. Categorical \leftrightarrow Dirichlet
4. Normal \leftrightarrow Normal

Conjugate prior for (univariate) Gaussian with known variance :-

likelihood $P(D|\theta) \sim N(\mu, \sigma^2)$, prior $P(\theta) \sim N(\mu_0, \sigma_0^2)$

$D = \{x_1, \dots, x_n\}$

$$f(\mu) = \frac{1}{\sqrt{2\pi\sigma^2}} \exp\left\{-\frac{1}{2\sigma^2}(\mu - \mu_0)^2\right\}$$

$$P(\theta|D) \propto P(D|\theta) P(\theta) \propto \exp\left\{-\frac{1}{2\sigma^2} \sum_{i=1}^n (x_i - \mu)^2 - \frac{1}{2\sigma_0^2}(\mu - \mu_0)^2\right\}$$

what will be $\tilde{\sigma}^2, \tilde{\mu}$? $\propto \exp\left\{-\frac{1}{2\tilde{\sigma}^2} \sum (x_i - \tilde{\mu})^2\right\} \sim N(\tilde{\mu}, \tilde{\sigma}^2)$

MAP estimate for linear regression

$$y_i = \omega^T x_i + \epsilon_i \sim N(\omega^T x_i, \sigma^2) \quad \text{where, } \epsilon_i \sim N(0, \sigma^2)$$

$$P(\theta|D) \propto \exp\left\{-\frac{1}{2\sigma^2} \sum (y_i - \omega^T x_i)^2\right\}$$

$$\hat{\theta}_{\text{MLE}} = \arg \min_{\omega} \sum_{i=1}^n (y_i - \omega^T x_i)^2$$

Now, $P(\omega) \sim N(0, \frac{1}{\lambda} I)$
 λ : hyperparameter.

$$P(\omega) = \frac{1}{(2\pi/\lambda)^{d/2}} \exp\left\{-\frac{1}{2} \omega^T \omega\right\}$$

$$\propto \exp\left\{-\frac{\lambda}{2} \|\omega\|^2\right\}$$

Now, $P(\omega|D) \propto P(D|\omega) P(\omega)$

Multivariate Gaussian distribution:

For $\bar{x} \in \mathbb{R}^d$, covariance matrix $\Sigma \in \mathbb{R}^{d \times d}$

$$P(\bar{x}) \sim N(\bar{\mu}, \Sigma)$$

where, $\sum_{i,j} \Sigma_{ij} = \sigma_{ij} = \text{cov}(x_i, x_j)$

$$P(x) = \frac{1}{(2\pi)^{d/2} |\Sigma|^{1/2}} \exp\left\{-\frac{1}{2} (\bar{x} - \bar{\mu})^T \Sigma^{-1} (\bar{x} - \bar{\mu})\right\}$$

$$\arg \max_{\omega} \left(\log P(D|\omega) + \log P(\omega) \right) :$$

$$\arg \min_{\omega} \left\{ \underbrace{\frac{1}{2\sigma^2} \sum_{i=1}^n (y_i - \omega^\top x_i)^2}_{\text{deterministic linear regression}} + \underbrace{\frac{\lambda}{2} \|\omega\|^2}_{\text{Regularizer.}} \right\} \equiv \left\{ \frac{1}{2\sigma^2} \|x\omega - y\|^2 + \frac{\lambda}{2} \|\omega\|^2 \right\}$$