



LDA 过程：假设给定数据：

$$f_k(x) = \frac{1}{(2\pi)^{p/2} |\Sigma_k|^{1/2}} \exp\left(-\frac{1}{2}(x - \mu_k)^T \Sigma_k^{-1} (x - \mu_k)\right)$$

- Assume the classes share a common covariance  $\Sigma_k = \Sigma, \forall k$

Compare two classes  $k$  and  $l$

$$\log \frac{Pr(G=k|X=x)}{Pr(G=l|X=x)} = \log \frac{f_k(x)}{f_l(x)} + \log \frac{\pi_k}{\pi_l}$$

$$= \log \frac{\pi_k}{\pi_l} - \frac{1}{2}(\mu_k - \mu_l)^T \Sigma^{-1} (\mu_k - \mu_l) + x^T \Sigma^{-1} (\mu_k - \mu_l)$$

二次项由于共同的协方差而消失

#### Parameter estimation

$\hat{\pi}_k = N_k/N$ , where  $N_k$  is the number of class- $k$  observations:

$$\hat{\mu}_k = \sum_{g_i=k} x_i / N_k;$$

$$\hat{\Sigma} = \sum_{k=1}^K \sum_{g_i=k} (x_i - \hat{\mu}_k)(x_i - \hat{\mu}_k)^T / (N - K)$$

$$\hat{\mu}_k = \frac{N_k}{N} \sum_{g_i=k} x_i$$

$$\hat{\Sigma} = \frac{N_k}{N(N-1)} \sum_{g_i=k} (x_i - \hat{\mu}_k)(x_i - \hat{\mu}_k)^T$$

Pooled covariance (合并方差)

$$\hat{\Sigma} = \frac{(N_1 - 1)\hat{\Sigma}_1 + (N_2 - 1)\hat{\Sigma}_2 + \dots + (N_K - 1)\hat{\Sigma}_K}{(N_1 - 1) + (N_2 - 1) + \dots + (N_K - 1)}$$

Weighted average

#### QDA: Quadratic Discriminant Analysis

Assume that each class has a specific covariance  $\Sigma_k$

Discriminant Functions:

$$\delta_k(x) = -\frac{1}{2} \log |\Sigma_k| - \frac{1}{2}(x - \mu_k)^T \Sigma_k^{-1} (x - \mu_k) + \log \pi_k$$

#### Probability and Estimator

贝叶斯网络

##### MLE: Maximum Likelihood Estimate

Choose  $\theta$  that maximizes probability of observed data  $D$

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

##### MAP: Maximum a Posterior

Choose  $\theta$  that is most probable given prior probability and the data

$$\text{MAP: } P(\theta|D) = \frac{P(D|\theta)P(\theta)}{P(D)} = \dots$$

CS182 CheatSheet

⇒ 边界:  $\log \frac{P(G=1|X=x)}{P(G=0|X=x)}$

$$= \frac{\frac{1}{2}}{\frac{1}{2}} - \frac{1}{2}(\hat{\mu}_1 + \hat{\mu}_2) \sum_{i=1}^2 (x_i - \hat{\mu}_i)^T$$

$$= 0.18 \left[ -\sum_{i=1}^2 (x_i - \hat{\mu}_i)^T \right] = 0$$

$$\Rightarrow \{x_1, x_2\} \mid x_1 + x_2 = 0.75 \}$$

Maximum Likelihood Estimate



$$x_1 \sim \theta, x_2 \sim 1-\theta$$

$$P(X=x) = \theta^x (1-\theta)^{1-x}$$

Maximum A Posteriori (MAP) Estimate



Data set  $D$  of independent, identically distributed (iid) flips:  $\alpha_1$  ones,  $\alpha_0$  zeros

$$P(D|\theta) = P(\alpha_1, \alpha_0|\theta) = \theta^{\alpha_1} (1-\theta)^{\alpha_0}$$

Assume prior  $P(\theta) = \text{Beta}(\beta_1, \beta_0) = \frac{\theta^{\beta_1-1} (1-\theta)^{\beta_0-1}}{\Gamma(\beta_1+\beta_0)}$

Then

$$\hat{\theta}^{MAP} = \arg \max_{\theta} P(D|\theta)P(\theta) = \frac{\alpha_1 + \beta_1 - 1}{(\alpha_1 + \beta_1 - 1) + (\alpha_0 + \beta_0 - 1)}$$

(like MLE, but halving  $\beta_1-1$  additional heads,  $\beta_0-1$  additional tails)

#### Naive Bayes

Train Naive Bayes:

◦ for each value  $y_k$  3.3.假设: 所有变量均独立.

◦ estimate  $\pi_k = P(Y=y_k)$

◦ for each value  $x_{ij}$  of each attribute  $X_i$

◦ estimate  $\theta_{ijk} = P(X_i = x_{ij}|Y=y_k)$

Classify

◦  $Y^{\text{new}} \leftarrow \arg \max_{y_k} P(Y=y_k) \prod_i P(X_i^{\text{new}}|Y=y_k)$

◦  $Y^{\text{new}} \leftarrow \arg \max_{y_k} \pi_k \prod_i \theta_{ijk}$

#### Estimate Parameters

Maximum likelihood estimates:

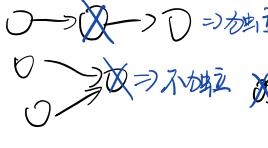
$$\hat{\pi}_k = \hat{P}(Y=y_k) = \frac{\#D(Y=y_k)}{|D|}$$

$$\hat{\theta}_{ijk} = \hat{P}(X_i = x_{ij}|Y=y_k) = \frac{\#D(X_i = x_{ij} \wedge Y=y_k)}{\#D(Y=y_k)}$$

MAP estimates (Beta, Dirichlet priors):

$$\hat{\pi}_k = \hat{P}(Y=y_k) = \frac{\#D(Y=y_k) + (\beta_k - 1)}{|D| + \sum_m (\beta_m - 1)}$$

$$\hat{\theta}_{ijk} = \hat{P}(X_i = x_{ij}|Y=y_k) = \frac{\#D(X_i = x_{ij} \wedge Y=y_k) + (\beta_k - 1)}{\#D(Y=y_k) + \sum_m (\beta_m - 1)}$$



CS182 CheatSheet

6

#### 高斯贝叶斯

##### Continuous $X_i$ (Gaussian Naive Bayes)

Assume:

$$P(X_i = x|Y=y_k) = \frac{1}{\sigma_{ik}\sqrt{2\pi}} \exp\left(-\frac{(x - \mu_{ik})^2}{2\sigma_{ik}^2}\right)$$

$$\text{Maximum likelihood estimates: } \begin{aligned} \hat{\mu}_{ik} &= \frac{1}{\sum_j \delta(Y^j = y_k)} \sum_j X_j^j \delta(Y^j = y_k) \\ \text{jth feature} &\quad \text{kth class} \\ \text{ith feature} &\quad \text{jth training example} \end{aligned}$$

$$\hat{\sigma}_{ik}^2 = \frac{1}{\sum_j \delta(Y^j = y_k)} \sum_j (X_j^j - \hat{\mu}_{ik})^2 \delta(Y^j = y_k)$$

感知机 perceptron

#### Kernel Method & SVM 感知机

The online Learning Model: Perceptron 对于每个数据点，如果没出错则  $w_i^{th} = w_i^t$

Algorithm

• Set  $t = 1$ , start with all zero vector  $w_t$

• Given example  $x$ , predict positive iff  $w_t \cdot x \geq 0$

• On a mistake, update:

- Mistake on positive, then update  $w_{t+1} \leftarrow w_t + x$
- Mistake on negative, then update  $w_{t+1} \leftarrow w_t - x$

• Note:  $w_t = a_{11}x_{11} + \dots + a_{ik}x_{ik}$

而  $w^t$  是一个矩阵  $\Rightarrow w_t = a_{11}x_{11} + a_{12}x_{12} + \dots + a_{ik}x_{ik}$

$$\text{E.g. } x[-1, 2, -1, 0, 1, 1, 0, -1, 1, -1]$$

$$y = - + + - - - +$$

$$\textcircled{1} \text{ sign}(w_1^t x_1) = 0 \quad \textcircled{2} \text{ sign}(w_2^t x_2) > 0 \quad \textcircled{3} \text{ sign}(w_3^t x_3) < 0 \quad \textcircled{4} \text{ sign}(w_4^t x_4) < 0 \quad \textcircled{5} \text{ sign}(w_5^t x_5) < 0 \quad \textcircled{6} \text{ sign}(w_6^t x_6) < 0 \quad \textcircled{7} \text{ sign}(w_7^t x_7) < 0 \quad \textcircled{8} \text{ sign}(w_8^t x_8) < 0 \quad \textcircled{9} \text{ sign}(w_9^t x_9) < 0$$

不等

- Definition: The margin of example  $x$  w.r.t. a linear sep.  $w$  is the distance from  $x$  to the plane  $w \cdot x = 0$
- The margin  $\gamma_w$  of a set of examples  $S$  w.r.t. a linear separator  $w$  is the smallest margin over points  $x \in S$ .
- The margin  $\gamma$  of  $n$  examples  $S$  is the maximum  $\gamma_w$  over all linear seps  $w$

#### Mistake Bound

Theorem: If data linearly separable by margin  $\gamma$  and points inside a ball of radius  $R$ , then Perceptron makes  $R/\gamma^2$  mistakes.

#### Proof:

Idea: analyze  $w_t \cdot w^*$  and  $\|w_t\|$ , where  $w^*$  is the max-margin sep,  $\|w^*\| = 1$ .

Claim 1:  $w_{t+1} \cdot w^* \geq w_t \cdot w^* + \gamma$ . (because  $l(x)x \cdot w^* \geq \gamma$ )

Claim 2:  $\|w_{t+1}\|^2 \leq \|w_t\|^2 + R^2$ . (by Pythagorean Theorem)

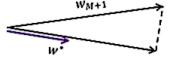
After  $M$  mistakes:

$w_{M+1} \cdot w^* \geq \gamma M$  (by Claim 1)

$\|w_{M+1}\| \leq R\sqrt{M}$  (by Claim 2)

$w_{M+1} \cdot w^* \leq \|w_{M+1}\|$  (since  $w^*$  is unit length)

$$\text{So, } \gamma M \leq R\sqrt{M}, \text{ so } M \leq \left(\frac{R}{\gamma}\right)^2.$$



#### SVM

Directly optimize for the maximum margin separator:

- Input:  $S = \{(x_1, y_1), \dots, (x_m, y_m)\}$

• Maximize  $\gamma$  under:

◦  $\|w\|^2 = 1$

◦  $\forall i, y_i w \cdot x_i \geq \gamma$

i.e.  $w^* = \frac{w}{\gamma}$ , 把优化问题转化成凸问题：

• Maximize  $\|w\|^2$  under  $\forall i, y_i w \cdot x_i \geq 1$

• If data isn't perfectly linearly separable?

◦ Replace "# mistakes" with upper bound called "hinge loss"

◦ Minimize  $\|w\|^2 + C \sum_i \xi_i$  s.t.  $\forall i, y_i w \cdot x_i \geq 1 - \xi_i, \xi_i \geq 0$

#### Lagrangian Dual of SVMs

$$\begin{aligned} W &= (0, 0) \\ \textcircled{1} w_2 &= w_1 - (-1, 2) = (1, -1) \\ \textcircled{2} w_3 &= w_1 + (1, 1) = (2, 1) \\ \textcircled{3} w_4 &= w_1 - (-1, -1) = (2, 1) \end{aligned}$$

### Solution

We first formulate the Lagrangian function of the primal problem:

$$L(\mathbf{w}, \xi, \alpha, \lambda) = \|\mathbf{w}\|_2^2 - \sum_{i=1}^n \alpha_i(y_i \mathbf{w}^\top \mathbf{x}_i - 1),$$

where  $\alpha_i \geq 0$  ( $\forall i$ ) is the dual variable. Because strong duality holds in the primal problem, the optimal optimization variables  $\{\mathbf{w}^*, \alpha^*\}$  should satisfy K.K.T. conditions:

- primal:  $y_i \mathbf{w}^\top \mathbf{x}_i \geq 1, \forall i$ ,
- dual:  $\alpha_i^* \geq 0, \forall i$ ,
- complementary:  $\alpha_i^*(y_i \mathbf{w}^\top \mathbf{x}_i - 1) = 0, \forall i$ ,
- stationary:  $\nabla_{\mathbf{w}^*} L = 0$ .

According to the stationary condition, we have

$$\nabla_{\mathbf{w}} L = 2\mathbf{w} - \sum_{i=1}^n \alpha_i y_i \mathbf{x}_i = 0, \Rightarrow \mathbf{w} = \frac{1}{2} \sum_{i=1}^n \alpha_i y_i \mathbf{x}_i,$$

Substituting them into the Lagrangian function yields the dual function  $g(\alpha, \lambda)$ ,

$$\begin{aligned} g(\alpha) &= \inf_{\mathbf{w}, \xi} L(\mathbf{w}, \alpha) \\ &= \frac{1}{4} \left( \sum_{i=1}^n \alpha_i y_i \mathbf{x}_i \cdot \sum_{j=1}^n \alpha_j y_j \mathbf{x}_j \right) - \sum_{i=1}^n \alpha_i \left( y_i \left( \frac{1}{2} \sum_{j=1}^n \alpha_j y_j \mathbf{x}_j \cdot \mathbf{x}_i \right) - 1 \right) \\ &= -\frac{1}{2} \sum_{i=1}^n \sum_{j=1}^n y_i y_j \alpha_i \alpha_j \langle \mathbf{x}_i, \mathbf{x}_j \rangle + \sum_{i=1}^n \alpha_i. \end{aligned}$$

Thus, the dual problem is

$$\begin{aligned} \max_{\alpha} & -\frac{1}{2} \sum_{i=1}^n \sum_{j=1}^n y_i y_j \alpha_i \alpha_j \langle \mathbf{x}_i, \mathbf{x}_j \rangle + \sum_{i=1}^n \alpha_i, \\ \text{s.t. } & \alpha_i \geq 0, \forall i. \end{aligned}$$

### Kernels

$K(\cdot, \cdot)$  is a kernel if it can be viewed as a legal definition of inner product:

$$\exists \phi : X \rightarrow \mathbb{R}^N, \text{ s.t. } K(x, z) = \phi(x) \cdot \phi(z)$$

### Theorem (Mercer)

$K$  is a kernel if and only if:

- $K$  is symmetric
- For any set of training points  $x_1, x_2, \dots, x_m$  and for any  $a_1, a_2, \dots, a_m \in \mathbb{R}$ , we have:

$$\sum_{i,j} a_i a_j K(x_i, x_j) \geq 0$$

$$a^T K a \geq 0$$

i.e.,  $K = (K(x_i, x_j))_{i,j=1,\dots,n}$  is positive semi-definite

### Kernels

- Linear:  $K(x, z) = x \cdot z$
- Polynomial:  $K(x, z) = (x \cdot z)^2$  or  $K(x, z) = (1 + x \cdot z)^d$
- Gaussian:  $K(x, z) = \exp \left[ -\frac{\|x-z\|^2}{2\sigma^2} \right]$
- Laplace:  $K(x, z) = \exp \left[ -\frac{\|x-z\|}{2\sigma^2} \right]$

## 神经网络:

首先画图, 找出传播链 反向传播四公式

然后根据链推导关系

找出所有路径, 算导数

## MLP 全连接 NN

BP1

If  $Y = \sigma(X)$ , then:

$$\frac{\partial L}{\partial X} = \frac{\partial L}{\partial Y} \odot \sigma'(x)$$

Where  $\odot$  means element-wise product

### BP2 & 3 & 4

If  $Y = W \cdot X + B$

$$\frac{\partial L}{\partial X} = W^T \cdot \frac{\partial L}{\partial Y}$$

$$\frac{\partial L}{\partial B} = \frac{\partial L}{\partial Y}$$

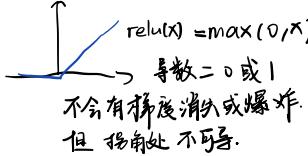
$$\frac{\partial L}{\partial W} = \frac{\partial L}{\partial Y} \cdot X^T$$

### 推导的核心

考虑:  $y_i = w_{i1}x_1 + w_{i2}x_2 + \dots + w_{in}x_n + b_i$

$$\frac{\partial y_i}{\partial w_{ij}} = x_j, \quad \frac{\partial y_i}{\partial x_j} = w_{ij}, \quad \frac{\partial y_i}{\partial b_i} = 1$$

还有一个 ReLU ← 常用损失函数的导数



### Softmax

If  $X = [x_1, \dots, x_n], Y = \text{SoftMax}(X) = [y_1, \dots, y_n]$

Softmax with 交义熵 :

$$\text{Thus } y_i = \frac{e^{x_i}}{\sum_{j=1}^n e^{x_j}} \text{ and obviously } \sum_{i=1}^n y_i = 1$$

We have:

$$\frac{\partial L}{\partial X} = Y - \hat{Y}$$

我们有:

$$\frac{\partial y_i}{\partial x_j} = \begin{cases} y_i(1 - y_i), & i = j \\ -y_i y_j, & i \neq j \end{cases}$$

### Tanh

$$\text{We have } \tau(x) = \tanh(x) = \frac{e^x - e^{-x}}{e^x + e^{-x}}$$

It's equal to:

$$\tau'(x) = 1 - \tau^2(x)$$

$$\frac{\partial Y}{\partial X} = \text{diag}(Y) - Y^T Y$$

### Sigmoid

We have  $\sigma(x) = \text{sigmoid}(x) = \frac{1}{1+e^{-x}}$

$$\sigma'(x) = \sigma(x)(1 - \sigma(x))$$

## PCA 阵容

### Dimension Reduction: PCA

Denoted that  $v_1, \dots, v_d$  are the d principal components,  $v_i \cdot v_j = \begin{cases} 1, & i = j \\ 0, & i \neq j \end{cases}$

Let  $X = [x_1, \dots, x_n]$  (Columns are the datapoints)

Maximizes sample variance of projected data

$$\frac{1}{n} \sum_{i=1}^n (\mathbf{v}^\top \mathbf{x}_i)^2 = \mathbf{v}^\top \mathbf{X} \mathbf{X}^\top \mathbf{v} \text{ s.t. } \mathbf{v}^\top \mathbf{v} = 1$$

Lagrangian:  $\arg \max_{\mathbf{v}} \mathbf{v}^\top \mathbf{X} \mathbf{X}^\top \mathbf{v} - \lambda \mathbf{v}^\top \mathbf{v}$

We finally have  $(\mathbf{X} \mathbf{X}^\top - \lambda I)\mathbf{v} = 0$ .  $(\mathbf{X} \mathbf{X}^\top)\mathbf{v} = \lambda \mathbf{v}$

**Maximum Variance Direction:** 1<sup>st</sup> PC a vector  $v$  such that projection on to this vector capture maximum variance in the data (out of all possible one dimensional projections)

$$\frac{1}{n} \sum_{i=1}^n (\mathbf{v}^\top \mathbf{x}_i)^2 = \mathbf{v}^\top \mathbf{X} \mathbf{X}^\top \mathbf{v}$$

**Minimum Reconstruction Error:** 1<sup>st</sup> PC a vector  $v$  such that projection on to this vector yields minimum MSE reconstruction

$$\text{blue}^2 + \text{green}^2 = \text{black}^2$$

black<sup>2</sup> is fixed (it's just the data)

So, maximizing blue<sup>2</sup> is equivalent to minimizing green<sup>2</sup>

### PCA 过程:

① 计算 center design matrix:  $\bar{\mathbf{x}}$   
计算  $\bar{\mathbf{x}}$ , 对每一个  $\mathbf{x}_i - \bar{\mathbf{x}}$ , 组成  $\bar{\mathbf{x}}$ :  $\begin{bmatrix} \mathbf{x}_1 - \bar{\mathbf{x}} \\ \mathbf{x}_2 - \bar{\mathbf{x}} \\ \vdots \\ \mathbf{x}_n - \bar{\mathbf{x}} \end{bmatrix}$

② 计算  $\bar{\mathbf{x}}^\top \bar{\mathbf{x}}$

③  $\det(sI - \bar{\mathbf{x}}^\top \bar{\mathbf{x}}) = 0$  解出 S 的取值

即为方差最大方向

而方差为对应的 S 值

这里的  $\mathbf{v}^\top \mathbf{X}$  是投影, 再乘上  $\mathbf{v}$  就是在投影上的向量。

- The eigenvalue  $\lambda$  denotes the amount of variability captured along that dimension.
- Zero eigenvalues indicate no variability along those directions  $\Rightarrow$  data lies exactly on a linear subspace
- Only keep data projections onto principal components with non-zero eigenvalues, say  $v_1, \dots, v_k$ , where  $k = \text{rank}(\mathbf{X} \mathbf{X}^\top)$

### Clustering: KMeans (KMeans 聚类)

### Denotions

Given a sample  $\mathcal{X} = \{\mathbf{x}^t\}_{t=1}^N$

Find  $k$  reference vectors  $\mathbf{m}_j$  which best represent the data

### Encoding/Decoding

Each data point  $\mathbf{x}^t$  is represented by the index  $i$  of the nearest reference vector  $i = \arg \min_j \|\mathbf{x}^t - \mathbf{m}_j\|$

We can use labels  $\mathbf{b}^t$  for  $\mathbf{x}^t$  as:

$$b_i^t = \begin{cases} 1 & \text{if } i = \arg \min_j \|\mathbf{x}^t - \mathbf{m}_j\| \\ 0 & \text{otherwise} \end{cases}$$

The Total Reconstruction Error:

$$E(\{\mathbf{n}_i\}_{i=1}^k | \mathcal{X}) = \sum_t \sum_i b_i^t \|\mathbf{x}^t - \mathbf{m}_i\|^2$$

$$\arg \min_{\{\mathbf{m}_i\}_{i=1}^k, \{\mathbf{b}^t\}_{t=1}^N} \sum_t \sum_i b_i^t \|\mathbf{x}^t - \mathbf{m}_i\|^2$$

However, since  $b_i^t$  also depends on  $\mathbf{m}_i$ , the optimization problem cannot be solved analytically, but iteratively

② 选迭代中心点。

① 根据距离最近原则分类。

② 计算每个类的中间点 (均值)

③ 重复第一步, 迭代

定义 Loss =  $\frac{1}{N} \sum_i \|\mathbf{x}^i - \mathbf{m}_{c(i)}\|^2$

这里指分类后距对应类别中心点的距离

### Algorithm

**Initialize**  $\mathbf{m}_i, i = 1, \dots, k$  (e.g.,  $k$  randomly selected  $\mathbf{x}^t$ )

**Repeat**

For all  $\mathbf{x}^t \in \mathcal{X}$ , we obtain the estimated labels

$$b_i^t = \begin{cases} 1 & \text{if } i = \arg \min_j \|\mathbf{x}^t - \mathbf{m}_j\| \\ 0 & \text{otherwise} \end{cases}$$

For all  $\mathbf{m}_i, i = 1, \dots, k$ , we obtain (by taking the derivative of  $E(\{\mathbf{m}_i\}_{i=1}^k | \mathcal{X})$  with respect to  $\mathbf{m}_i$  and setting it to 0)

$$\mathbf{m}_i = \frac{\sum_t b_i^t \mathbf{x}^t}{\sum_t b_i^t}$$

The reference vector is set to the mean (center) of all the instances that it represents.

**Until**  $\mathbf{m}_i$  converge.

► Note  $\mathbf{m}_i$  here is the same as the formula for the mean estimation in classification, except that we place the estimated labels  $\mathbf{b}^t$  in place of the labels  $\mathbf{r}^t$ .

## 梯度下降 Gradient Descent

### Basic Problem

$$\arg \min_{x \in \mathbb{R}^n} f(x) \quad \text{每层大小: } W = \frac{W \times F + 2P}{S} + 1 \quad H = \frac{H \times F + 2P}{S} + 1 \Rightarrow W \times H \times K$$

Iteration:  $x^{r+1} = x^r - \gamma_r \cdot \nabla f(x^r)$

### Convexity

$$\begin{aligned} f(\lambda x) + (1 - \lambda)y &\leq \lambda f(x) + (1 - \lambda)f(y) \\ f(x) &\geq f(y) + \nabla f(y)^T(x - y) \\ \nabla^2 f(x) &\geq 0 \end{aligned}$$

每层参数:  $(F \times F \times D + 1) \times K$

全连接层:  $(\text{Feature\_in} + 1) \times \text{Feature\_out}$ .

### L-smooth

$$\|\nabla f(x) - \nabla f(y)\| \leq L \|x - y\|$$

### Descent Lemma

$$|f(x) - f(y) - \nabla f(y)^T(x - y)| \leq \frac{L}{2} \|x - y\|^2$$

If  $f$  twice-differentiable, L-smooth  $\Leftrightarrow \nabla^2 f(x) \leq LI$ ,  $\mathbf{d}^T \nabla^2 f(x) \mathbf{d} \leq L \|\mathbf{d}\|^2, \forall \mathbf{d}$

### Convergence Analysis

Optimality measure  $M(x^r)$

- Convex:  $\|x^r - x^*\|$  or  $f(x^r) - f^*$

- Non-convex:  $\|\nabla f(x^r)\|$

Order of convergence  $q$  s.t.

$$\sup \left\{ q \mid \lim_{r \rightarrow +\infty} \frac{M(x^{r+1})}{M(x^r)^q} < \infty \right\}$$

- $q = 1$ : Linear convergence,  $q = 2$ : quadratic

Rate of Convergence: Given  $q$ ,

$$\lim_{r \rightarrow \infty} \frac{M(x^{r+1})}{M(x^r)^q} = n$$

- Sublinear: Rate = 1, Superlinear: Rate = 0

### Convergence under Convexity

#### Theorem

Let  $f$  be convex with bounded gradient  $\|\nabla f(x^r)\| \leq B$ , then the sequence  $(x^r)_{r \in \mathbb{N}}$  generated by GD with step size  $\gamma_r = \frac{\sqrt{B}}{\|\nabla f(x^r)\|}$  satisfies

$$\min_{r=0, \dots, T-1} f(x^r) - f^* \leq \frac{B \|x^0 - x^*\|}{\sqrt{T}}.$$

### Convergence under Smoothness

#### Theorem

Let  $f$  be  $L$ -smooth, then the sequence  $(x^r)_{r \in \mathbb{N}}$  generated by GD with step size  $\gamma \leq 1/L$  satisfies

$$\min_{r=0, \dots, T-1} \|\nabla f(x^r)\|^2 \leq \frac{\frac{2}{\gamma} (f(x^0) - f^*)}{T}.$$

### Convexity & Smoothness

#### Theorem

Let  $f$  be convex and  $L$ -smooth, then the sequence  $(x^r)_{r \in \mathbb{N}}$  generated by GD with step size  $\gamma \leq 1/L$  satisfies

$$f(x^T) - f^* \leq \frac{\|x^0 - x^*\|}{2\gamma T}.$$

### Upper & Lower Bound

### $\mu$ -Strong Convexity

$$f(x) \geq f(y) + \nabla f(y)^T(x - y) + \frac{\mu}{2} \|x - y\|^2$$

### L-Smoothness

$$f(x) \leq f(y) + \nabla f(y)^T(x - y) + \frac{L}{2} \|x - y\|^2$$

### Implication

$$\begin{aligned} \nabla f(x^r)^T(x^r - x^*) &\geq f(x^r) - f_{\frac{L}{2}}(x^r - x^*) \\ f(x^{r+1}) &\leq f(x^r) - \frac{r}{2} \|\nabla f(x^r)\|^2 \end{aligned}$$

## 拉格朗日

### Basics

#### Original Function

Minimize  $f_0(x)$  with optimal  $p^*$  subject to  $f_i(x) \leq 0, i = 1, \dots, m$

We have  $L(x, \lambda) = f_0(x) + \lambda_1 f_1(x) + \dots + \lambda_m f_m(x)$

#### Dual Function

•

Lower Bound Property: if  $\lambda \geq 0$  and  $x$  are primal feasible:

$$g(\lambda) \leq f_0(x)$$

#### Dual Problem

Maximize  $g(\lambda)$  with optimal  $d^*$  subject to  $\lambda \geq 0$

We have:  $d^* \leq p^*$ , denoted that  $p^* - d^*$  the optimal dual gap.

The Convex Problem have  $p^* = d^*$

#### KKT Optimal Condition

|  |                 |
|--|-----------------|
| $f_i(x^*) \leq 0$  | Primal Feasible |
| $\lambda_i^* f_i(x^*) = 0$                               | Dual Feasible   |
| $\nabla f_0(x^*) + \sum \lambda_i^* \nabla f_i(x^*) = 0$ | Complementary   |
|  | Stationary      |

#### Equality Constraints

$$\begin{aligned} \partial A &= \emptyset \\ \partial(\alpha X) &= \alpha \partial X \\ \partial(X + Y) &= \partial X + \partial Y \\ \partial(\text{Tr}(X)) &= \text{Tr}(\partial X) \\ \partial(XY) &= (\partial X)Y + X\partial Y \\ \partial(X^{-1}) &= -X^{-1}(\partial X)X^{-1} \\ \partial(\det(X)) &= \text{Tr}(\text{adj}(X) \partial X) \\ &= \det(X) \text{Tr}(X^{-1} \partial X) \\ \text{In } \det(X) &= \det(X) \text{Tr}(X^{-1} \partial X) \\ \partial(X^T) &= (\partial X)^T \end{aligned}$$

$$\frac{\partial X^T a}{\partial X} = \frac{\partial a^T X}{\partial X} = a$$

$$\frac{\partial a^T X b}{\partial X} = ab^T$$

$$\frac{\partial a^T X^T b}{\partial X} = ba^T$$

$$\frac{\partial X}{\partial x_{ij}} = J_{ij}$$

Minimize  $f_0(x)$  subject to  $f_i(x) \leq 0, i = 1, \dots, m$  and  $h_i(x) = 0, i = 1, \dots, p$

We have  $L(x, \lambda, v) = f_0(x) + \sum \lambda_i f_i(x) + \sum v_i h_i(x)$

• Dual Function:  $g(\lambda, v) = \inf_x L(x, \lambda, v)$

• Dual Problem: Maximize  $g(\lambda, v)$  subject to  $\lambda \geq 0$

KKT:

|   |                 |
|---|-----------------|
| $f_i(x^*) \leq 0, h_i(x^*) = 0$   | Primal Feasible |
| $\lambda_i^* \geq 0$  | Dual Feasible   |
| $\lambda_i^* f_i(x^*) = 0$  | Complementary   |
| $\nabla f_0(x^*) + \sum \lambda_i^* \nabla f_i(x^*) + \sum v_i^* \nabla h_i(x^*) = 0$ | Stationary      |

**LSTM:**

