

Concentration Inequalities

Markov's inequality

For RV X with finite $\mathbb{E}[X]$ and $X \geq 0$:

$$\mathbb{P}(X \geq t) \leq \mathbb{E}[X]/t \quad t > 0$$

Pf. Refer to the property

$$\begin{aligned} \int_0^\infty \mathbb{P}(X \geq x) dx &= x \mathbb{P}(X \geq x) \Big|_0^\infty - \int_0^\infty x d\mathbb{P}(X \geq x) \\ &= 0 + \int_0^\infty f_X(x) dx = \mathbb{E}[X] \end{aligned}$$

From here, consider

$$\int_0^\infty \mathbb{P}(X \geq x) dx \geq \int_0^t \mathbb{P}(X \geq t) dx \geq t \mathbb{P}(X \geq t)$$

Chebyshev's inequality

For RV X with finite $\mathbb{E}[X]$, $\text{Var}(X)$ and $X \geq 0$:

$$\mathbb{P}(|X - \mathbb{E}[X]| \geq t) \leq \text{Var}(X)/t^2 \quad t > 0$$

Pf. Apply Markov's with $Y = |X - \mathbb{E}[X]|^2$

Chernoff's bound

For RV X with finite MGF = $\mathbb{E}[e^{\lambda X}]$, for any $t \geq 0$:

$$\mathbb{P}(X \geq t) = \mathbb{P}(e^{\lambda X} \geq e^{\lambda t}) \leq \mathbb{E}[e^{\lambda X - \lambda t}]$$

$$\mathbb{P}(X \leq -t) = \mathbb{P}(e^{-\lambda X} \geq e^{\lambda t}) \leq \mathbb{E}[e^{-\lambda X - \lambda t}]$$

Pf. Apply Markov's with $Y = e^{\pm \lambda X}$ (notice $Y \geq 0$). One may select a λ that to maximize $\mathbb{E}[e^{\lambda X - \lambda t}]$.

Chernoff's bound corollaries

- C^2 -Subgaussian RV $\mathbb{E}[X] = 0$ and $\mathbb{E}[e^{\lambda X}] \geq e^{\lambda^2 C^2/2}$, using Chernoff's bound: $\mathbb{P}(X \geq t) \leq e^{\lambda^2 C^2/2 - \lambda t}$. Setting $\lambda = \frac{t}{C^2}$ yields the tightest bound: $\mathbb{P}(|X| \geq t) \leq 2e^{-t^2/(2C^2)}$
- If $X \in [a, b]$, then $X - \mathbb{E}[X]$ is $\frac{(a-b)^2}{2}$ subgaussian.

Hoeffding's bound

For IID RV $\in [a, b]$ with mean μ and empirical mean $\hat{\mu}_n$,

$$\mathbb{P}(|\hat{\mu}_n - \mu| \geq t) \leq 2e^{\frac{-2nt^2}{(b-a)^2}} \quad t > 0$$

Pf. $\hat{\mu}_n - \mu$ is $\frac{(b-a)^2}{4n}$ subgaussian (use the corollary).

Estimation

Emperical Mean and Variances

Let $X \in [[a, b]$ be a random variable, $\bar{\mu}$ be the empirical mean, $\hat{\sigma}^2$ the empirical variance computed with μ , and $\hat{\sigma}^2$ with $\bar{\mu}$ respectively. One may bound $|\bar{\mu} - \mu|$ with the Hoeffding's bound on X_i , $|\hat{\sigma}^2 - \sigma^2|$ with the Hoeffding's bound on $Y_i = |X_i - \mu|^2$, and $|\hat{\sigma}^2 - \sigma^2|$ using the identity $(\hat{\sigma}^2 - \sigma^2) = (\hat{\sigma}^2 - \sigma^2) + (\bar{\mu} - \mu)^2$

Optimization

Consider a function $f : \mathbb{R}^d \rightarrow \mathbb{R}$ which we want to optimize wrt to its input \mathbf{x} : $\hat{\mathbf{x}} = \min_{\mathbf{x} \in \mathbb{R}^d} f(\mathbf{x})$.

Derivative test

Gradient the first derivative $\nabla f(\mathbf{w})$ is defined by

$$\left[\frac{\partial f}{\partial w_1}(\mathbf{w}) \quad \frac{\partial f}{\partial w_2}(\mathbf{w}) \quad \cdots \quad \frac{\partial f}{\partial w_d}(\mathbf{w}) \right]^\top$$

Hessian the second derivative $\nabla^2 f(\mathbf{w})$ is defined by

$$\begin{bmatrix} \frac{\partial^2 f}{\partial w_1^2}(\mathbf{w}) & \frac{\partial^2 f}{\partial w_1 \partial w_2}(\mathbf{w}) & \cdots & \frac{\partial^2 f}{\partial w_1 \partial w_d}(\mathbf{w}) \\ \frac{\partial^2 f}{\partial w_2 \partial w_1}(\mathbf{w}) & \frac{\partial^2 f}{\partial w_2^2}(\mathbf{w}) & \cdots & \frac{\partial^2 f}{\partial w_2 \partial w_d}(\mathbf{w}) \\ \vdots & \vdots & \ddots & \vdots \\ \frac{\partial^2 f}{\partial w_d \partial w_1}(\mathbf{w}) & \frac{\partial^2 f}{\partial w_d \partial w_2}(\mathbf{w}) & \cdots & \frac{\partial^2 f}{\partial w_d^2}(\mathbf{w}) \end{bmatrix}$$

Local min: \mathbf{x}^* s.t. for some neighborhood around \mathbf{x}^* , $f(\mathbf{x}) \geq f(\mathbf{x}^*)$ for all \mathbf{x} from the neighborhood

Global min: \mathbf{x}^* s.t. $\forall \mathbf{x} \in \mathbb{R}^d, f(\mathbf{x}) \geq f(\mathbf{x}^*)$

Stationary points: \mathbf{x} s.t. $\nabla f(\mathbf{x}) = \mathbf{0}$. It can either be a local max, min, or saddle point.

Conditions for local min: On top of $\nabla f(\mathbf{x}) = \mathbf{0}$:

- (necessary) $\mathbf{v}^\top \nabla^2 f(\mathbf{x}) \mathbf{v} \geq 0, \forall \mathbf{v} \in \mathbb{R}^d$
- (sufficient) $\mathbf{v}^\top \nabla^2 f(\mathbf{x}) \mathbf{v} > 0, \forall \mathbf{v} \in \mathbb{R}^d - \{\mathbf{0}\}$

Alternatively, the entries of $\nabla^2 f(\mathbf{x})$ must be ≥ 0

Convexity

f is convex if for any $\mathbf{w}_1, \mathbf{w}_2 \in \mathbb{R}^d$ and $\alpha \in (0, 1)$

$$f(\alpha \mathbf{w}_1 + (1 - \alpha) \mathbf{w}_2) \leq \alpha f(\mathbf{w}_1) + (1 - \alpha) f(\mathbf{w}_2)$$

If f is differentiable, then

- Once: $f(\mathbf{w}_1) \geq f(\mathbf{w}_2) + \nabla f(\mathbf{w}_2)^\top (\mathbf{w}_1 - \mathbf{w}_2)$
- Twice: $\nabla^2 f(\mathbf{x})$ is positive semi-definite everywhere.

Jensen's inequality if f is convex, for any distribution D

$$f(\mathbb{E}_{\mathbf{w} \sim D}[\mathbf{w}]) \leq \mathbb{E}_{\mathbf{w} \sim D}[f(\mathbf{w})]$$

Lipchitz smooth $\|\nabla f(\mathbf{w}_1) - \nabla f(\mathbf{w}_2)\|_2 \leq L \|\mathbf{w}_1 - \mathbf{w}_2\|_2$. All L -smooth function must satisfy $f(\mathbf{w}_1) \leq f(\mathbf{w}_2) + \nabla f(\mathbf{w}_2)^\top (\mathbf{w}_1 - \mathbf{w}_2) + \frac{L}{2} \|\mathbf{w}_1 - \mathbf{w}_2\|_2^2$

Gradient Descent

Assume f is once differentiable, $\mathbf{w}_{t+1} \leftarrow \mathbf{w}_t + \eta \nabla f(\mathbf{w}_t)$.

Gradient Descent Convergence

If f if L -smooth, gradient descent converges $\propto \frac{1}{t}$:

$$\begin{aligned} f(\mathbf{w}') &\leq f(\mathbf{w}) + \nabla f(\mathbf{w})^\top (\mathbf{w}' - \mathbf{w}) + \frac{L}{2} \|\mathbf{w}' - \mathbf{w}\|_2^2 \\ \text{Substituting } \mathbf{w} &= \mathbf{w}_t \text{ and } \mathbf{w}' = \mathbf{w}_{t+1} = \mathbf{w}_t - \eta_t \nabla f(\mathbf{w}_t), \\ f(\mathbf{w}_{t+1}) &\leq f(\mathbf{w}_t) - \left(\eta_t - \frac{\eta_t^2 L}{2} \right) \|\nabla f(\mathbf{w}_t)\|_2^2 \end{aligned}$$

Choose $0 < \eta_t < \frac{2}{L}, f(\mathbf{w}_{t+1}) - f(\mathbf{w}_t) \leq \frac{-\eta}{2} \|\nabla f(\mathbf{w}_t)\|_2^2$
Telescopic sum and rearrangement:

$$\begin{aligned} \sum_{i=0}^t \|\nabla f(\mathbf{w}_i)\|_2^2 &\leq \frac{2}{\eta} (f(\mathbf{w}_0) - f(\mathbf{w}_{t+1})) \\ &\leq \frac{2}{\eta} (f(\mathbf{w}_0) - f_{\min}) \\ \min_i \|\nabla f(\mathbf{w}_i)\|_2^2 &\leq \frac{2(f(\mathbf{w}_0) - f_{\min})}{\eta \left(\frac{t+1}{2} (f(\mathbf{w}_0) - f_*) \right)} \end{aligned}$$

if we want gradient $\leq \epsilon$, set $t \geq \left\lceil \frac{2(f(\mathbf{w}_0) - f_*)}{\epsilon^2} \right\rceil$

Reducing the Optimality Gap

If f is also convex, the error converges $\propto \frac{1}{t}$:

$$f(\mathbf{w}_*) \geq f(\mathbf{w}_t) + \nabla f(\mathbf{w}_t)^\top (\mathbf{w}_* - \mathbf{w}_t).$$

Expanding the square:

$$\begin{aligned} \frac{1}{2} \|\mathbf{w}_{t+1} - \mathbf{w}_*\|_2^2 &= \frac{1}{2} \|\mathbf{w}_t - \mathbf{w}_*\|_2^2 \\ &\quad - \eta \nabla f(\mathbf{w}_t)^\top (\mathbf{w}_t - \mathbf{w}_*) + \frac{\eta^2}{2} \|\nabla f(\mathbf{w}_t)\|_2^2 \end{aligned}$$

Convexity implies

$$\eta \nabla f(\mathbf{w}_t)^\top (\mathbf{w}_* - \mathbf{w}_t) \leq \eta (f(\mathbf{w}_*) - f(\mathbf{w}_t))$$

and L -smoothness implies (proven prior)

$$\frac{\eta^2}{2} \|\nabla f(\mathbf{w}_t)\|_2^2 \leq \eta (f(\mathbf{w}_t) - f(\mathbf{w}_{t+1}))$$

put altogether

$$f(\mathbf{w}_{t+1}) - f(\mathbf{w}_*) \leq \frac{1}{2\eta} (\|\mathbf{w}_t - \mathbf{w}_*\|_2^2 - \|\mathbf{w}_{t+1} - \mathbf{w}_*\|_2^2)$$

Summing telescopically and using how $f(\mathbf{w}_t)$ decreases:

$$\begin{aligned} f(\mathbf{w}_{T+1}) - f(\mathbf{w}_*) &\leq \frac{1}{2\eta(T+1)} \|\mathbf{w}_0 - \mathbf{w}_*\|_2^2 \\ \text{if we want LHS} &\leq \epsilon, \text{ set } t = \left\lceil \frac{1}{2\eta\epsilon} \|\mathbf{w}_0 - \mathbf{w}_*\|_2^2 \right\rceil \end{aligned}$$

Stochastic Gradient Descent

Let $f(\mathbf{w}) = \mathbb{E}_{x,y \sim D}[\ell(\mathbf{w}; \mathbf{x}, \mathbf{y})]$ and the grad $\nabla_{\mathbf{w}} = \mathbb{E}_{x,y \sim D}[\nabla_{\mathbf{w}} \ell(\mathbf{w}; \mathbf{x}, \mathbf{y})]$. Define $\mathbf{g}_{\mathbf{w}} = \nabla_{\mathbf{w}} \ell(\mathbf{w}; \mathbf{x}, \mathbf{y})$. Then, $\nabla_{\mathbf{w}} = \mathbb{E}[\mathbf{g}_{\mathbf{w}}]$. Finally, force $\mathbb{E}[\|\mathbf{g}_{\mathbf{w}}\|_2^2] \leq G^2, \forall \mathbf{w}$:

$$\begin{aligned} \mathbb{E}_t[\|\mathbf{w}_{t+1} - \mathbf{w}_*\|_2^2] &= \mathbb{E}_t\left[\frac{1}{2} \|\mathbf{w}_t - \mathbf{w}_*\|_2^2 - \eta \mathbf{g}_t^\top (\mathbf{w}_t - \mathbf{w}_*) + \frac{\eta^2}{2} \|\mathbf{g}_t\|_2^2\right] \\ &= \frac{1}{2} \mathbb{E}_t[\|\mathbf{w}_t - \mathbf{w}_*\|_2^2] - \eta \mathbb{E}_t[\mathbf{g}_t^\top (\mathbf{w}_t - \mathbf{w}_*)] + \frac{\eta^2}{2} \mathbb{E}_t[\|\mathbf{g}_t\|_2^2] \\ &\leq \frac{1}{2} \mathbb{E}_t[\|\mathbf{w}_t - \mathbf{w}_*\|_2^2] - \eta \nabla_{\mathbf{w}} f(\mathbf{w})^\top (\mathbf{w}_t - \mathbf{w}_*) + \frac{\eta^2}{2} G^2 \\ \eta \nabla_{\mathbf{w}} f(\mathbf{w})^\top (\mathbf{w}_t - \mathbf{w}_*) &\leq \frac{1}{2} \mathbb{E}_t[\|\mathbf{w}_t - \mathbf{w}_*\|_2^2] - \frac{1}{2} \mathbb{E}_t[\|\mathbf{w}_{t+1} - \mathbf{w}_*\|_2^2] + \frac{\eta^2}{2} G^2 \\ \eta (f(\mathbf{w}_t) - f(\mathbf{w}_*)) &\leq (\text{convexity of } f) \\ \eta \mathbb{E}_t[(f(\mathbf{w}_t) - f(\mathbf{w}_*))] &\leq \frac{1}{2} \mathbb{E}_t[\|\mathbf{w}_t - \mathbf{w}_*\|_2^2] - \frac{1}{2} \mathbb{E}_t[\|\mathbf{w}_{t+1} - \mathbf{w}_*\|_2^2] + \frac{\eta^2}{2} G^2 \end{aligned}$$

Summing telescopically and divide both sides by $\eta(T+1)$:

$$\begin{aligned} \frac{1}{T+1} \sum_{t=0}^T \mathbb{E}[f(\mathbf{w}_t) - f(\mathbf{w}_*)] &\leq \frac{\|\mathbf{w}_0 - \mathbf{w}_*\|_2^2}{2\eta(T+1)} + \frac{\eta G^2}{2} \\ \text{Let } \overline{\mathbf{w}}_T &= \frac{1}{T+1} \sum_{t=0}^T \mathbf{w}_t \text{ and } \eta = \frac{\|\mathbf{w}_0 - \mathbf{w}_*\|_2}{G\sqrt{T+1}}, \text{ we pick} \\ T &= \left\lceil \frac{\|\mathbf{w}_0 - \mathbf{w}_*\|_2^2 + G^2}{2\epsilon^2} \right\rceil \text{ if we want error } \leq \epsilon \end{aligned}$$

PAC Learning

Define error: $\text{err}_{\mathcal{D}}(h) := \mathbb{P}_{\mathbf{x} \sim \mathcal{D}}(c(\mathbf{x}) \neq h(\mathbf{x}))$.

Learning Finite Classes

For finite \mathcal{C} , the probability error $\geq \epsilon$ yet h is chosen from n data points is $(1 - \epsilon)^n$. Suppose this for some h is $|\mathcal{C}|(1 - \epsilon)^n \leq \delta$ gives $n = \left\lceil \frac{\log(|\mathcal{C}|) + \log(1/\delta)}{\epsilon} \right\rceil$.

Bounding Empirical Error

Define $\hat{\text{err}}_{\mathcal{D}}(h) := \sum_{i=1}^n \mathbb{I}\{c(\mathbf{x}_i) \neq h(\mathbf{x}_i)\}$. Hoeffding's bound give, for an h , $\mathbb{P}(|\hat{\text{err}}(h) - \text{err}(h)| \geq \epsilon) \leq 2e^{-2n\epsilon^2}$. Union bound for all classes gives $2|\mathcal{C}|e^{-2n\epsilon^2} \leq \delta$ which gives $n = \left\lceil \frac{1}{2\epsilon^2} \log\left(\frac{2|\mathcal{C}|}{\delta}\right) \right\rceil$.

Empirical Error Trick

Given $\mathbb{P}(\forall h, |\hat{\text{err}}(h) - \text{err}(h)| \leq \epsilon/2) \geq 1 - \delta$. Let \hat{h} optimizes $\hat{\text{err}}$ and h_* optimizes $\text{err}(h)$. This means $\text{err}(\hat{h}) \leq \hat{\text{err}}(\hat{h}) + \epsilon/2 \leq \hat{\text{err}}(h_*) + \epsilon/2 \leq \text{err}(h_*) + \epsilon$

Perceptron Algorithm

Assuming $y_i = \pm 1$ and $\|\mathbf{x}_i\|_2^2, \|\mathbf{w}_*\|_2^2 \leq 1$,

$$\exists \mathbf{w}_*, y_i = \text{sign}(\mathbf{w}_*^\top \mathbf{x}_i) \text{ and } |\mathbf{w}_*^\top \mathbf{x}| \geq \gamma$$

The perceptron algorithm starts with $\mathbf{w}_0 = \mathbf{0}$ and updates $\mathbf{w}_{t+1} \leftarrow \mathbf{w}_t + y_i \mathbf{x}_i$ for a misclassified i . First, $\mathbf{w}_*^\top \mathbf{w}_{t+1} = \mathbf{w}_*^\top (\mathbf{w}_t + y_i \mathbf{x}_i) = \mathbf{w}_*^\top \mathbf{w}_t + y_i \mathbf{w}_*^\top \mathbf{x}_i \geq \mathbf{w}_*^\top \mathbf{w}_t + \gamma$ so $\mathbf{w}_*^\top \mathbf{w}_t \geq t\gamma$. Moreover, $\|\mathbf{w}_{t+1}\|_2^2 = \|\mathbf{w}_t + y_i \mathbf{x}_i\|_2^2 = \|\mathbf{w}_t\|_2^2 + 2y_i \mathbf{w}_t^\top \mathbf{x}_i + y_i^2 \|\mathbf{x}_i\|_2^2$. Since $y_i \mathbf{w}_t^\top \mathbf{x}_i \leq 0, \|\mathbf{w}_{t+1}\|_2^2 \leq \|\mathbf{w}_t\|_2^2 + 1$ so $\|\mathbf{w}_t\|_2 \leq \sqrt{t}$. Put together, $t \leq \frac{1}{\gamma^2}$ is when the model converges.

Linear Regression

Suppose $\mathbf{X} \in \mathbb{R}^{n \times d}$ and $\mathbf{y} \in \mathbb{R}^{n \times 1}$, there exists $\mathbf{w}_* \in \mathbb{R}^{d \times 1}$ such that $\mathbf{y} = \mathbf{X} \mathbf{w}_* + \boldsymbol{\eta}$ for $\boldsymbol{\eta} \sim \mathcal{N}(\mathbf{0}, \sigma^2 \mathbf{I})$.

MLE Solution (OLS)

The MLE gives the training objective to be $\ell(\mathbf{w}) = \|\mathbf{y} - \mathbf{X} \mathbf{w}\|_2^2$. The gradient is $2(\mathbf{X}^\top \mathbf{X} \mathbf{w} - \mathbf{X}^\top \mathbf{y})$ and hessian $2\mathbf{X}^\top \mathbf{X}$, giving $\mathbf{w} = (\mathbf{X}^\top \mathbf{X})^{-1} \mathbf{X}^\top \mathbf{y}$.

Sample Complexity

$$\begin{aligned} \|\mathbf{y} - \mathbf{X} \hat{\mathbf{w}}\|_2^2 &= \|\mathbf{X} \hat{\mathbf{w}} - \mathbf{X} \mathbf{w}_*\|_2^2 \\ &\quad - 2\eta^\top \mathbf{X} (\hat{\mathbf{w}} - \mathbf{w}_*) + \|\boldsymbol{\eta}\|_2^2 \end{aligned}$$

$$\|\mathbf{y} - \mathbf{X} \hat{\mathbf{w}}\|_2^2 \leq \|\boldsymbol{\eta}\|_2^2 \text{ because of MLE}$$

$$\|\mathbf{X} (\hat{\mathbf{w}} - \mathbf{w}_*)\|_2^2 \leq 2\eta^\top \mathbf{X} (\hat{\mathbf{w}} - \mathbf{w}_*)$$

$$\|\mathbf{X} (\hat{\mathbf{w}} - \mathbf{w}_*)\|_2 \leq \frac{2\eta^\top \mathbf{X} (\hat{\mathbf{w}} - \mathbf{w}_*)}{\|\mathbf{X} (\hat{\mathbf{w}} - \mathbf{w}_*)\|_2}$$

Let $r = \text{rank}(\mathbf{X}) = \text{rank}(\mathbf{X}^\top \mathbf{X})$, then there exists orthonormal bases $\boldsymbol{\Phi} \in \mathbb{R}^{n \times r}$ such that $\boldsymbol{\Phi}^\top \boldsymbol{\Phi} = \mathbf{I}_r$ and $\mathbf{X} (\hat{\mathbf{w}} - \mathbf{w}_*) = \boldsymbol{\Phi} \mathbf{z}$ for some $\mathbf{z} \in \mathbb{R}^r$. Cauchy-Schwartz:

$$\frac{\eta^\top \mathbf{X} (\hat{\mathbf{w}} - \mathbf{w}_*)}{\|\mathbf{X} (\hat{\mathbf{w}} - \mathbf{w}_*)\|_2} = \frac{\eta^\top \boldsymbol{\Phi} \mathbf{z}}{\|\boldsymbol{\Phi} \mathbf{z}\|_2} = \frac{\eta_r^\top \mathbf{z}}{\|\mathbf{z}\|_2} \leq \|\boldsymbol{\eta}_r\|_2$$

Since η_r is normal, $\mathbb{E}[\|\eta_r\|] = \sigma\sqrt{r}$, implying $\|\mathbf{X}(\hat{\mathbf{w}} - \mathbf{w}_*)\|_2^2 \leq 4\|\eta_r\|_2^2$ and $\mathbb{E}[\|\mathbf{X}(\hat{\mathbf{w}} - \mathbf{w}_*)\|_2^2] \leq 4\sigma^2 r$. Define $V := \frac{1}{n} \|\mathbf{X}\hat{\mathbf{w}} - \mathbf{X}\mathbf{w}_*\|_2^2 - 4\frac{\sigma^2 r}{n} \leq \frac{4}{n} \|\eta_r\|_2^2 - 4\frac{\sigma^2 r}{n}$, then $\mathbb{E}\left[e^{\lambda V}\right] \leq \left(\mathbb{E}\left[e^{\frac{4\lambda}{n}(\eta_r^2/\sigma^2 - 1)}\right]\right)^r$. $\eta_{r,1}^2/\sigma^2$ is standard normal. We use the MGF for $z^2 - 1$ if $z \in \mathcal{N}(0, 1)$ as follows: $\mathbb{E}[e^{\lambda(z^2 - 1)}] \leq e^{2\lambda}$, set $\lambda = \frac{4\lambda\sigma^2}{n}$, $\mathbb{E}[e^{\lambda V}] \leq e^{\frac{32r\lambda^2\sigma^4}{n^2}}$.

Agnostic Linear Regression Model

Suppose \mathbf{w}_* leads to the smallest true error $\text{OPT} := \text{err}(\mathbf{w}_*)$. We first bound $\text{err}(\hat{\mathbf{w}}) \leq 4$ because of normalization constraint. For a given \mathbf{w} , Hoeffding's give

$$\mathbb{P}(|\hat{\text{err}}(\mathbf{w}) - \text{err}(\mathbf{w})| \geq t) \leq 2 \exp\left(\frac{-2nt^2}{4}\right)$$

Quantization define a grid $\mathcal{H}_{\epsilon'}$ of size ϵ' , $|\mathcal{H}_{\epsilon'}| = (\frac{2}{\epsilon'})^d$.

$$\mathbb{P}(\exists \mathbf{w} \in \mathcal{H}_{\epsilon'} : |\hat{\text{err}}(\mathbf{w}) - \text{err}(\mathbf{w})| \geq t) \leq 2\left(\frac{2}{\epsilon'}\right)^d e^{(-\frac{2nt^2}{4})}$$

Set $t = \epsilon$ and RHS $\leq \delta$ gives $n = O(\frac{d \ln(1/\epsilon') + \ln(1/\delta)}{\epsilon^2})$. So $\text{err}(\hat{\mathbf{v}}) \leq \text{err}(\mathbf{v}_*) + \epsilon/2$ when \mathbf{v} nearest to \mathbf{w} in $\mathcal{H}_{\epsilon'}$.

$$\begin{aligned} \text{err}(\mathbf{v}) &= \mathbb{E}_{(\mathbf{x}, y)} \left[(y - \mathbf{v}^\top \mathbf{x})^2 \right] \\ &= \mathbb{E}_{(\mathbf{x}, y)} \left[(y - \mathbf{w}^\top \mathbf{x} + \mathbf{w}^\top \mathbf{x} - \mathbf{v}^\top \mathbf{x})^2 \right] \\ &= \mathbb{E}_{(\mathbf{x}, y)} \left[(y - \mathbf{w}^\top \mathbf{x})^2 \right] + \mathbb{E}_{(\mathbf{x}, y)} \left[(\mathbf{w}^\top \mathbf{x} - \mathbf{v}^\top \mathbf{x})^2 \right] \\ &\quad + 2\mathbb{E}_{(\mathbf{x}, y)} \left[(y - \mathbf{w}^\top \mathbf{x})(\mathbf{w}^\top \mathbf{x} - \mathbf{v}^\top \mathbf{x}) \right] \\ &\leq \text{err}(\mathbf{w}) + \mathbb{E}_{(\mathbf{x}, y) \sim \mathcal{D}} \left[(\mathbf{w}^\top \mathbf{x} - \mathbf{v}^\top \mathbf{x})^2 \right] \\ &\quad + 2\sqrt{\text{err}(\mathbf{w}) \cdot \mathbb{E}_{(\mathbf{x}, y) \sim \mathcal{D}} [(\mathbf{w}^\top \mathbf{x} - \mathbf{v}^\top \mathbf{x})^2]} \\ &= \text{err}(\mathbf{w}) + \text{err}_q + 2\sqrt{\text{err}(\mathbf{w}) \cdot \text{err}_q} \\ &\leq \text{err}(\mathbf{w}) + \text{err}_q + 4\sqrt{\text{err}_q} \end{aligned}$$

Therefore, we want $\text{err}_q + 4\sqrt{\text{err}_q} \leq \epsilon/4$:

$$\begin{aligned} \text{err}_q &= \mathbb{E}_{(\mathbf{x}, y)} \left[(\mathbf{w}^\top \mathbf{x} - \mathbf{v}^\top \mathbf{x})^2 \right] \\ &= \mathbb{E}_{(\mathbf{x}, y)} \left[\|\mathbf{w} - \mathbf{v}\|_2^2 \|\mathbf{x}\|_2^2 \right] \\ &\leq \|\mathbf{w} - \mathbf{v}\|_2^2 \cdot \mathbb{E}_{(\mathbf{x}, y)} \left[\|\mathbf{x}\|_2^2 \right] \\ &\leq \left(\frac{\epsilon}{2}\right)^2 \cdot \|\mathbf{1}_d\|_2^2 = \frac{(\epsilon)^2 d}{4} \end{aligned}$$

As such, $\epsilon' \in O(d/\sqrt{\epsilon})$, $n = O(\frac{d \ln(d/\epsilon) + \ln(1/\delta)}{\epsilon^2})$

Total Variation Distance

Given observations from $\mathbf{p} = (p_1, \dots, p_k)$, we measure the closeness of an estimate $\hat{\mathbf{p}} = (\hat{p}_1, \dots, \hat{p}_k)$:

$$d_{\text{TV}}(\hat{\mathbf{p}}, \mathbf{p}) = \max_{S \subseteq \{1, 2, \dots, k\}} \sum_{i \in S} (p_i - \hat{p}_i)$$

Break it down to positive and negative terms:

$$\sum_{i: p_i > \hat{p}_i} |p_i - \hat{p}_i| - \sum_{i: p_i < \hat{p}_i} |\hat{p}_i - p_i|$$

This value is maximal when $p_i > \hat{p}_i$. Since positive and negative sums are equal in magnitude,

$$\max_{S \subseteq \{1, 2, \dots, k\}} \sum_{i \in S} (p_i - \hat{p}_i) = \frac{1}{2} \sum_{i=1}^k |p_i - \hat{p}_i|$$

Estimator: $X_i : i = 1, 2, \dots, n$, we assign a Bernoulli random variable for each class where $\hat{p}_i = \frac{\sum_{j=1}^n \mathbb{I}(X_j = i)}{n}$.

$$\begin{aligned} d_{\text{TV}}(\hat{\mathbf{p}}, \mathbf{p}) &= \max_{S \subseteq \{1, 2, \dots, k\}} \sum_{i \in S} (p_i - \hat{p}_i) \\ &= \max_{S \subseteq \{1, 2, \dots, k\}} \sum_{i \in S} \frac{\mathbb{I}(X_j = i)}{n} - \sum_{i \in S} p_i \end{aligned}$$

By linearity of expectation, $\mathbb{E}\left[\sum_{i \in S} \frac{\mathbb{I}(X_j = i)}{n}\right] = \sum_{i \in S} p_i$.

This means $\sum_{i \in S} \frac{\mathbb{I}(X_j = i)}{n}$ is a random variable with mean $\sum_{i \in S} p_i$ and range $[0, 1]$. Applying Hoeffding's bound,

$$\mathbb{P}\left(\left|\sum_{i \in S} \frac{\mathbb{I}(X_j = i)}{n} - \sum_{i \in S} p_i\right| \geq \epsilon\right) \leq 2e^{-2\epsilon^2 n}$$

$$\mathbb{P}(|d_{\text{TV}}(\hat{\mathbf{p}}, \mathbf{p})| \geq \epsilon) \leq 2^k 2e^{-2\epsilon^2 n}$$

2^k is the number of possible S . Hence, $2^k 2e^{-2\epsilon^2 n} \leq \delta$ yields the bound for $n \geq \frac{k \log(2) + \log(2\delta)}{2\epsilon^2}$

Hypothesis Testing

H_0 : null hypothesis (nothing happens) H_1 : alternative. Define $\alpha = \mathbb{P}(\text{reject } H_0 | H_0)$ (significance) and $1 - \beta = \mathbb{P}(\text{reject } H_1 | H_0)$ (power). p -value: $p = \mathbb{P}(\text{event} | H_0 \text{ is accepted})$. $p \leq \alpha$ means we reject H_0 .

Unknown Null Hypothesis

Given two unknown distr p_* (reference) and q_* (test). $q_1, q_2, \dots, q_n \sim q_*$ Let $s = \frac{1}{n} \sum_{i=1}^n d_{\text{TV}}(p_*, q_i)$ and $\bar{s} = \mathbb{E}_{q_*}[s]$. Hoeffding gives

$$p = \mathbb{P}(|s - \bar{s}| \geq |s_{\text{obs}} - \bar{s}|) \leq 2e^{-2n|s_{\text{obs}} - \bar{s}|^2}$$

if accept $H_0 := \{p_* = q_*\}$, sample $p_1, p_2, \dots, p_N \sim p_*$, define $\hat{s} = \frac{1}{n} \sum_{i=1}^n d_{\text{TV}}(p_*, p_i)$: $\mathbb{P}(|\hat{s} - \bar{s}| \geq \epsilon) \leq 2e^{-2N\epsilon^2}$

$$\begin{aligned} p &\leq \mathbb{P}(|s - \bar{s}| \geq |s_{\text{obs}} - \hat{s}| - |\hat{s} - \bar{s}|) \\ &\leq \mathbb{P}(|s - \bar{s}| \geq |s_{\text{obs}} - \hat{s}| - |\hat{s} - \bar{s}| \mid |\hat{s} - \bar{s}| \leq \epsilon) \\ &\quad + \mathbb{P}(|\hat{s} - \bar{s}| \geq \epsilon) \text{ (total probability)} \\ &\leq \mathbb{P}(|s - \bar{s}| \geq |s_{\text{obs}} - \hat{s}| - \epsilon) + \mathbb{P}(|\hat{s} - \bar{s}| \geq \epsilon) \\ &\leq 2e^{2n(|s_{\text{obs}} - \bar{s}| - \epsilon)^2} + 2e^{2N\epsilon^2} \end{aligned}$$

choose $\epsilon = \frac{|s_{\text{obs}} - \bar{s}|}{2}$ and $n = N$: $p \leq 4e^{\frac{n|s_{\text{obs}} - \bar{s}|^2}{2}}$

Multiple Null Hypotheses

Given $H_0^{(i)}, i = 1, 2, \dots, K$, suppose each has significance α . Define FWER := $\mathbb{P}(\geq 1 \text{ false rejections}) \leq \sum_{i=1}^K \mathbb{P}(\text{reject } H_0^{(i)} | H_0^{(i)}) = \alpha K$, we set $\alpha \rightarrow \alpha/K$

Improving Success Probability

Testing algorithm \mathcal{A} passes with prob $1 - q$, $q \in [0, 0.5]$, using n samples (depends on k and ϵ). Suppose

$\mathcal{A}_i \sim \text{Ber}(q)$ and $\mathcal{A}' = \mathbb{I}\{\sum_{i=1}^m \mathcal{A}_i \geq \frac{m}{2}\}$.

$$\begin{aligned} \mathbb{P}(\mathcal{A}' = 0) &= \mathbb{P}\left(\sum_{i=1}^m \mathcal{A}_i \geq \frac{m}{2}\right) \\ &= \mathbb{P}\left(\frac{1}{m} \sum_{i=1}^m \mathcal{A}_i - q \geq \frac{1}{2} - q\right) \\ &\leq e^{-2m(\frac{1}{2} - q)^2} = e^{-\frac{m}{2}(1 - 2q)^2} \end{aligned}$$

Estimation Uniformity Testing

Let $\mathbf{u}_k = (1/k, 1/k, \dots, 1/k)$. Let our estimator be $\hat{\mathbf{p}}$ such that $d_{\text{TV}}(\hat{\mathbf{p}}, \mathbf{p}) \leq \epsilon/3$,

- If $\mathbf{p} = \mathbf{u}_k$, $d_{\text{TV}}(\hat{\mathbf{p}}, \mathbf{u}_k) \leq \epsilon/3$
- If $d_{\text{TV}}(\mathbf{p}, \mathbf{u}_k) \geq \epsilon$, d_{TV} is a metric so $d_{\text{TV}}(\hat{\mathbf{p}}, \mathbf{u}_k) \geq d_{\text{TV}}(\mathbf{p}, \mathbf{u}_k) - d_{\text{TV}}(\mathbf{p}, \hat{\mathbf{p}}) \geq 2\epsilon/3$. $\epsilon/2$ could be used.

Sample complexity: $n \geq \frac{k \log(2) + \log(2\delta)}{2\epsilon^2}$

Collision Uniformity Testing

Let $C = \sum_{i < j} \mathbb{I}\{X_i = X_j\}$ be the number of pairwise collisions among n samples from \mathbf{p} . $\mathbb{E}[C] = \binom{n}{2} \|\mathbf{p}\|_2^2$.

- If $\mathbf{p} = \mathbf{u}_k$, $\mathbb{E}[C] = \binom{n}{2} \cdot \frac{1}{k}$
- $d_{\text{TV}}(\mathbf{p}, \mathbf{u}_k) \geq \epsilon$, $\|\mathbf{p}\|_2^2 > \frac{1 + \epsilon^2}{k} \Rightarrow \mathbb{E}[C] > \binom{n}{2} \cdot \frac{1 + \epsilon^2}{k}$
- Set $T = \binom{n}{2} \cdot \frac{1 + \epsilon^2/2}{k}$, accept if $C < T$, reject if $C \geq T$

Claim: $\text{Var}(C) \leq n^2 \|\mathbf{p}\|_2^2 + n^3 \|\mathbf{p}\|_3^3$, $\|\mathbf{p}\|_3^3 = \sum_{i=1}^k p_i^3$. Observe that because $n \geq \sqrt{k}$, $\|\mathbf{p}\|_2 \geq \frac{1}{\sqrt{k}}$, and $\|\mathbf{p}\|_3 \leq \|\mathbf{p}\|_2$, we can further bound the variance by $\text{Var}[C] \leq 2n^3 \|\mathbf{p}\|_2^3$.

Using $\mathbb{P}(|C - \mathbb{E}[C]| \geq 2\sigma) \leq \frac{1}{4}$ and $\sigma = \sqrt{\text{Var}[C]} \leq \sqrt{2n^3 \|\mathbf{p}\|_2^3}$ we prove

- **Case 1 (Uniform):** $\mathbf{p} = \mathbf{u}_k$

$$\begin{aligned} \mathbb{E}[C] &= \binom{n}{2} \cdot \frac{1}{k}, \\ \sigma &\leq \sqrt{2n^3} \cdot \frac{1}{k^{3/2}} \end{aligned}$$

We want $C < T = \binom{n}{2} \cdot \frac{1 + \epsilon^2/2}{k}$, so:

$$\mathbb{E}[C] + 2\sigma < T \Rightarrow n = O\left(\frac{\sqrt{k}}{\epsilon^2}\right)$$

- **Case 2 (Far):** $d_{\text{TV}}(\mathbf{p}, \mathbf{u}_k) \geq \epsilon/2$

$$\|\mathbf{p}\|_2^2 \geq \frac{1 + \epsilon^2 + \alpha}{k},$$

$$\mathbb{E}[C] \geq \binom{n}{2} \cdot \frac{1 + \epsilon^2 + \alpha}{k}$$

We want $C \geq T$, so:

$$\mathbb{E}[C] - 2\sigma \geq T \Rightarrow n = O\left(\frac{\sqrt{k}}{\epsilon^2}\right)$$

Conclusion: $n = \Theta\left(\frac{\sqrt{k}}{\epsilon^2}\right)$ suffice to separate error $\leq 1/4$.

ℓ_2 Identity Testing

Goal: Given i.i.d. samples from unknown \mathbf{q} and explicit \mathbf{p} , test:

- H_0 : $\mathbf{p} = \mathbf{q}$ vs. H_1 : $\|\mathbf{p} - \mathbf{q}\|_1 \geq \epsilon$

Key idea: Use a test statistic that estimates $\|\mathbf{p} - \mathbf{q}\|_2^2$.

Tester: Let $X_i \sim \text{Pois}(nq_i)$, define

$$Z = \sum_{i=1}^k (X_i - np_i)^2 - X_i$$

Then:

$$\mathbb{E}[Z] = n^2 \|\mathbf{p} - \mathbf{q}\|_2^2 \quad (\text{bias-corrected estimator})$$

Theorem: With $n = \Theta\left(\frac{\sqrt{k}}{\epsilon^2}\right)$ samples, this tester distinguishes H_0 from H_1 with constant probability.

Why it works:

- Under H_0 , $\mathbb{E}[Z] = 0$
- Under H_1 , $\|\mathbf{p} - \mathbf{q}\|_1 \geq \epsilon \Rightarrow \|\cdot\|_2^2 \geq \epsilon^2/k$
- So $\mathbb{E}[Z] \geq n^2 \cdot \epsilon^2/k$
- With variance $\text{Var}(Z) = O(n^2)$, Chebyshev implies signal is detectable with $n = \Theta(\sqrt{k}/\epsilon^2)$