# Approximate Counting

Given $x_1, ..., x_n$, such that $x_i \in \Sigma = \{$ red, blue, ... $\}$

Count frequencies of red, blue, green
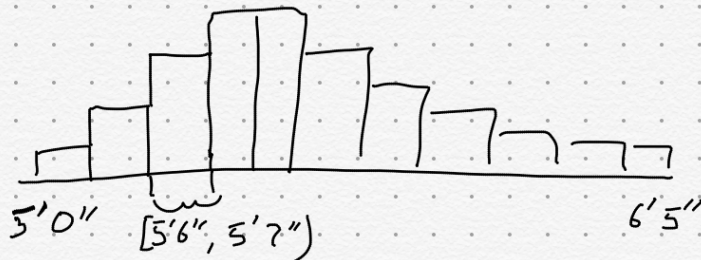- query $(\sigma) \to$ # of times $\sigma \in \Sigma$ shows up
- $O(|\Sigma|)$ for 100% acc
- Approximate with count-min sketch $(O(m) \quad m \ll |\Sigma|)$

What if $\Sigma \to \mathbb{Z}_{(0,D)}$ ?
- better to use histogram because of the <u>smoothness structure</u>
  - values near each other have similar # of occurances

Ex. Heights of students



$5'0"$ $[5'6", 5'7")$ $6'5"$

Query $(5'6")$ ⎫
Query $(5'7")$ ⎬ relatively closer
Query $(6'0")$ ⎭

# Aside: Over- and Under- fitting

Overfitting
- memorizing dataset / learning concepts that won't generalize
- ✱ Occurs when lots of free parameters

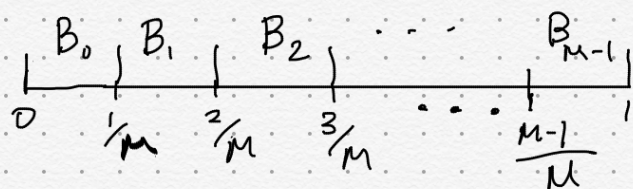Underfitting $\to$ Opposite
- too few parameters
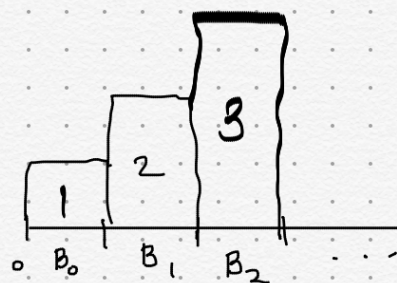
# Histograms

① Observe $x_1, \ldots, x_n$

   ↳ $x_i$'s are iid

   ↳ $x_i \in [0,1]$ (normalized data)

② "Sketch": Create histogram to $M$ bins



  ↳ maintain count over $n$ →

  ↳ Normalize counts to have a valid probability density

$$\hat{p}(B_\ell) = \left(\frac{M}{n}\right) |B_\ell|$$

     width of bucket

      ↳ count in bucket $\ell$

  ↳ $x_i$'s are sampled from density function $p(x)$

    ↳ $x_i \sim p(x)$

      ↳ $\hat{p}(x)$ approximates $p(x)$

        ↳ $O(M)$ parameters

        ↳ If $M$ is too small,

           ↳ Bad estimate because miss features of data (over smoothing)

        ↳ If $M$ is too big, bad estimate b/c risk noise

    ↳ <u>Intuition</u>

| | | Error |
|---|---|---|
| Bin is 100% of data | ⌊_____⌋ | $\dfrac{\sigma}{\sqrt{n}}$ |
| Bin is 1% of data | ⌊_⌋ | $\dfrac{\sigma}{\sqrt{1\% \text{ of } n}}$ |

# Histograms cont'd

→ $x_1, \ldots, x_n$ i.i.d where $x_i \sim p(x)$

→ Assume smoothness → <u>Lipschitz smoothness/continuity</u> → (L-Lipschitz)

$\quad\quad\quad\quad\quad\quad\quad\quad\quad\quad\quad\quad\quad\quad \hookrightarrow |p(x) - p(y)| \leq L|x-y|$

→ $\hat{p}(x) = \dfrac{M}{h} \cdot \underline{\text{count}(B_{(x)})}$

$\quad\quad\quad\quad\quad\quad\quad\quad \hookrightarrow$ count in bucket $x$ is in

<u>Goal</u>: compare $\hat{p}(x)$ to $p(x)$

$E(\hat{p}(x)) = M \cdot P(\underline{x \in B_\ell})$

$\quad\quad\quad\quad\quad\quad \hookrightarrow x$ in bucket $B_\ell$

$\quad\quad\quad = M \cdot \displaystyle\int_{\frac{\ell-1}{M}}^{\frac{\ell}{M}} p(v)\, dv$

$\quad\quad\quad = M \cdot \left[ F\left(\frac{\ell}{M}\right) - F\left(\frac{\ell-1}{M}\right) \right] \to F$ is cdf of $p(x)$

$\quad\quad\quad = \dfrac{F\left(\frac{\ell}{M}\right) - F\left(\frac{\ell-1}{M}\right)}{1/M} = \dfrac{F\left(\frac{\ell}{M}\right) - F\left(\frac{\ell-1}{M}\right)}{\frac{\ell}{M} - \frac{\ell-1}{M}} \left.\rule{0pt}{40pt}\right] \to$ looks like derivative

$\quad\quad\quad$ By the mean value thm

$\quad\quad\quad = p(x^*)$ for some $x^* \in \left[\frac{\ell-1}{M}, \frac{\ell}{M}\right)$

$E(\hat{p}(x)) = p(x^*)$

$E(\hat{p}(x)) - p(x) = p(x^*) - p(x) \quad\quad \to$ Bias related to smoothness and number of bins

$\quad\quad\quad\quad = L \cdot |x^* - x| \leq \dfrac{L}{M}$

# Histograms cont'd

Bias is $\frac{L}{M}$

Variance: $\text{Var}(\hat{p}(x)) = M^2 \cdot \text{var}\left(\frac{1}{n} \sum_{i=1}^{n} I(x_i \in B_\ell)\right)$

$$= M^2 \frac{P(x \in B_\ell) \cdot (1 - P(x \in B_\ell))}{n}$$

Assume roughly uniform distribution $\Rightarrow P(x \in B_\ell) = \frac{1}{M}$

$$\text{Var}(\hat{p}(x)) = M^2 \frac{\frac{1}{M}\left(1-\frac{1}{M}\right)}{n}$$

$$= \frac{M-1}{n}$$

$\text{MSE} = \text{Bias}^2 + \text{Var}$

$f(M) = \frac{L^2}{M^2} + \frac{M-1}{n}$  (for uniform distribution)

$f'(M) = \frac{-2L^2}{M^3} + \frac{1}{n} \implies$ Optimal $M$ is

$$M^3 = 2L^2 n$$

$$M = \sqrt[3]{2L^2 n} = O(\sqrt[3]{n})$$

For roughly uniform $p(x)$, $M_{opt} = \sqrt[3]{2L^2 n} = O(\sqrt[3]{n})$

Assuming nonuniformity,

$\text{Var} = M^2 \frac{P(x \in B_\ell)(1 - P(x \in B_\ell))}{n}$    $\max(P(x \in B_\ell)(1 - P(x \in B_\ell)) = \frac{1}{4}$

$$\text{Var} \le \frac{M^2}{4n}$$

$\therefore \text{MSE} = \frac{L^2}{M^2} + \frac{M^2}{4n}$

$f'(M) = \frac{-2L^2}{M} + \frac{2M}{4N}$

$M_{opt} = \sqrt[4]{4L^2 N} \implies O(M_{opt}) \le \frac{13}{10} \sqrt[3]{L^2 N}$ for roughly uniform    $\left(\approx \sqrt[3]{2}\right)$

$O(M_{opt}) \le \frac{3}{2} \sqrt[4]{L^2 N}$ for very not uniform    $\left(\approx \sqrt[4]{4}\right)$

Note: Smoothness term $L$ allows for $M_{opt}$ to be cubed/fourth root rather than linear with data. Intuition: know one point $\rightarrow$ know points nearby!!