# Any Single Gaussian

_Intro_

How to summarize data?

$x^1, x^2, \ldots \quad x^m.$

- mean $\mu = \frac{1}{m} \sum_i x^i$

- Variance $\sigma^2 = \frac{1}{m} \sum_i (x^i - \mu)^2$

What about in higher dim? mean is still average.

Variance?

Variance along direction (unit vector $v$):

Variance of $x^1 \cdot v, \; x^2 \cdot v, \; \ldots \; x^m \cdot v$

$$\sigma_v^2 = \frac{1}{m} \sum_i (x^i \cdot v - \mu \cdot v)^2$$

$$= \frac{1}{m} \sum_i \left((x^i - \mu) \cdot v\right)^2$$

$$= v^T \left(\frac{1}{m} \sum_i (x^i - \mu)(x^i - \mu)^T\right) v$$

$$\Sigma_{ii} = \mathbb{E}\left((X_i - \mu_i)^2\right)$$

$$\Sigma_{ij} = \mathbb{E}\left((X_i - \mu_i)(X_j - \mu_j)\right).$$

$\Sigma$ : covariance matrix

---

Probability distribution $\quad\quad Pr(X = a)$

Probability density function $\quad pdf(x = a)$
(continuous densities)

e.g. uniform in $[a, b]$ $\quad\quad p(x = t) = \begin{cases} \dfrac{1}{|b-a|} & t \in [a,b] \\ \\ 0 & t \notin [a,b] \end{cases}$

Can have $p(x) = f(x)$ where
$$f(x) \geq 0$$

nonnegative, integrable. $\quad\quad \int f < \infty.$

---

Most important example

| Gaussian Density. |

1-dim $\quad\quad -\dfrac{x^2}{2}$

$$\boxed{\text{Gaussian density.}} \quad \text{1-dim}$$

$$p(x) = \frac{1}{\sqrt{2\pi}} e^{-\frac{x^2}{2}}$$

Q. Why $\frac{1}{\sqrt{2\pi}}$ ?

— $x \in \mathbb{R}^d$ $\quad p(x) = \frac{1}{(\sqrt{2\pi})^d} e^{-\frac{\|x\|^2}{2}}$

$\|x\|^2 = \sum_i x_i^2$

$$= \prod_{i=1}^{d} \frac{1}{\sqrt{2\pi}} e^{-\frac{x_i^2}{2}}$$

Product of 1-dim Gaussians along each coordinate.

F1. for any unit vector $v$, $\quad x \sim N(0, I_d)$

$$X \cdot v \sim N(0, 1)$$

$$p(x \cdot v) = \int \frac{1}{(\sqrt{2\pi})^d} e^{-\frac{\|x\|^2}{2}} dx$$
$$=t$$

$$X : X \cdot v = t$$

Write $x$ in the basis $\{v, v_2, \dots v_d\}$

$$\|x\|^2 = \sum (x \cdot v_i)^2$$

$$\dots e^{-\frac{t}{2}} \int \quad 1 \quad e^{-\sum_{i=2}^{d} \frac{(x \cdot v_i)^2}{2}}$$

$$= \frac{1}{\sqrt{2\pi}} e^{-\frac{t^2}{2}} \int \frac{1}{(\sqrt{2\pi})} d_1 \, e^{-\sum_{i=2} \frac{()}{2}}$$

$$x \cdot v_1 = t$$

$$y = (x \cdot v_2 \ldots x \cdot v_{2t})$$
$$dy = |det[v_2 \ldots v_d]| \, dx$$

$$= \frac{e^{-\frac{t^2}{2}}}{\sqrt{2\pi}} \int \frac{1}{(\sqrt{2\pi})} d_1 \cdot e^{-\frac{\|y\|^2}{2}} dy$$

$$\underset{y}{}$$

$$= \frac{e^{-\frac{t^2}{2}}}{\sqrt{2\pi}} .$$

---

**F2.** $\quad \frac{1}{\sqrt{2\pi}} \int e^{-\frac{x^2}{2}} dx \quad = 1$

Let $\quad p(x)'' \quad p(x,y) = \frac{1}{2\pi} e^{-\frac{x^2+y^2}{2}}$

$$\frac{1}{2\pi} \int e^{-\frac{(x^2+y^2)}{2}} \, dxdy = \frac{1}{2\pi} \int_{r=0}^{\infty} 2\pi r \cdot e^{-\frac{r^2}{2}} dr$$

$$u = \frac{r^2}{2} \quad du = rdr$$

$$= \int_{u} e^{-u} du$$
$$= 1 .$$

$$\int p(x,y) = \left( \int p(x) \right)^2 \implies \int p(x) = 1 .$$

General Gaussian $\quad N(\mu, \sigma^2)$

$$p(x) = \frac{1}{\sqrt{2\pi}\,\sigma}\, e^{-\frac{(x-\mu)^2}{2\sigma^2}}$$

Multivariate $\quad N(\mu, \Sigma) = \dfrac{1}{(\sqrt{2\pi})^d \left|\det(\Sigma)\right|}\, e^{-\frac{(x-\mu)^T \Sigma^{-1}(x-\mu)}{2}}$

$$X \longrightarrow Ax + b$$

$$N(0, I_d) \longrightarrow N(b, AA^T)$$

$$X \sim N(\mu, \Sigma) \implies \Sigma^{-\frac{1}{2}}(X - \mu) \sim N(0, I_d)$$

---

P1. Estimate Gaussian given random iid samples.

$$X^1, X^2, \ldots\ldots X^m.$$

Natural choices: $\quad \tilde{\mu} = \dfrac{1}{m}\sum_i x^i$

$$\tilde{\Sigma} = \frac{1}{m} \sum (x^i - \tilde{\mu})(x^i - \tilde{\mu})^T$$

$$\tilde{\Sigma} = \frac{1}{m} \sum_i (x^i - \tilde{\mu})(x^i - \tilde{\mu})'$$

Something better?

**Goal**. Find $\tilde{\mu}, \tilde{\Sigma}$ s.t.

likelihood $N(\tilde{\mu}, \tilde{\Sigma})$ generates $x^1 ... x^m$

is maximized.

**Thm**. Empirical mean, covariance give

max likelihood estimates.

**Pf**. Consider the case when $\Sigma = \sigma^2 I$.

$$P(data) = \prod_{j=1}^m \frac{1}{(\sqrt{2\pi}\sigma)^d} e^{-\frac{\|x^j - \mu\|^2}{2\sigma^2}}$$

$$\log(P(data)) = -\log(\sqrt{2\pi}\sigma) - \sum_j \frac{\|x^j - \mu\|^2}{2\sigma^2}$$

$$\frac{\partial\downarrow}{\partial\sigma} = -\frac{md}{\sigma} + \frac{1}{\sigma^3} \sum_j \|x^j - \mu\|^2 = 0$$

$$- \text{''} \cdot \text{''} \cdot \|^2$$

$$\Rightarrow \quad \sigma^2 = \frac{1}{d} \cdot \frac{1}{m} \sum_j \| x^j - \mu \|^2$$

$$\frac{\partial \mathcal{L}}{\partial \mu_i} = 0 \quad \Rightarrow$$

$$-\frac{1}{2\sigma^2} \cdot 2 \sum_j (x_i^j - \mu_i) = 0$$

$$m \, \mu_i = \sum_i x_i^j$$

$$\mu_i = \frac{1}{m} \sum x_i^j$$

$$\mu = \frac{1}{m} \sum x^j .$$

---

**Thm.** [Central Limit Theorem] $X, X, \dots X$

iid      random variables from some distribution

with bounded variance    $\text{Var}(X_1) < \infty$.

Let $Y_n = \frac{1}{n} \sum_i X_i$

Then      $Y_n \longrightarrow N\left( \mathbb{E}(X_1), \frac{1}{n} \text{Var}(X_1) \right)$

Then $Y_n \longrightarrow N\left(E(X_1), \frac{1}{n}\text{Var}(X_1)\right)$

(converges in distribution)

---

This can be made more qualitative.

Thm. [Berry-Esseen] $X_1, \ldots X_n$ independ R.V.

$Y_n = \sum\limits_{i=1}^{n} X_i$ . $Z_n \sim N\left(E(Y_n), \text{Var}(Y_n)\right)$ .

Then, for any $t \in \mathbb{R}$,

$$\left| \Pr(Y_n \leq t) - \Pr(Z_n \leq t) \right| \leq C \cdot \frac{\sum\limits_{i=1}^{n} E(|X_i|^3)}{\text{Var}(Y_n)^{3/2}}$$

$C \in [0, 1]$ .

---

Example $X_i = \begin{cases} -1 & \text{w.p. } \frac{1}{2} \\ 1 & \text{w.p. } \frac{1}{2} \end{cases}$

$E(|X_i|^3) = 1$

$\text{Var}(Y_n) = n \, \text{Var}(X_1) = n \cdot \left(\frac{1}{2} \cdot 1 + \frac{1}{2} \cdot 1\right) = n$

$\therefore \forall t \qquad Z_n \sim (0, n)$

$$|P_1(Y_n \leq t) - P_1(Z_n \leq t)| \leq c \cdot \frac{n}{n^{3/2}} \leq \frac{c}{\sqrt{n}}.$$

"Invariance" Principle.

---

Note, bound also holds for $P_1(Y_n > t)$

$$|P_1(Y_n > t) - P_1(Z_n > t)| = |(1 - P_1(Y_n \leq t)) - (1 - P_1(Z_n \leq t))|$$

$$= |P_1(Y_n \leq t) - P_1(Z_n \leq t)|$$