

The Variance, the Standard Deviation and the Mean Square

Probability and Statistics for Data Science

Carlos Fernandez-Granda



These slides are based on the book [Probability and Statistics for Data Science](#) by Carlos Fernandez-Granda, available for purchase [here](#). A free preprint, videos, code, slides and solutions to exercises are available at <https://www.ps4ds.net>

Discrete random variable

The mean of a discrete random variable \tilde{a} with range A is

$$\mathbb{E} [\tilde{a}] := \sum_{a \in A} a p_{\tilde{a}}(a)$$

if the sum converges

Continuous random variable

The mean of a continuous random variable \tilde{a} is

$$\mathbb{E}[\tilde{a}] := \int_{a=-\infty}^{\infty} a f_{\tilde{a}}(a) \, da$$

if the integral converges

Sample mean

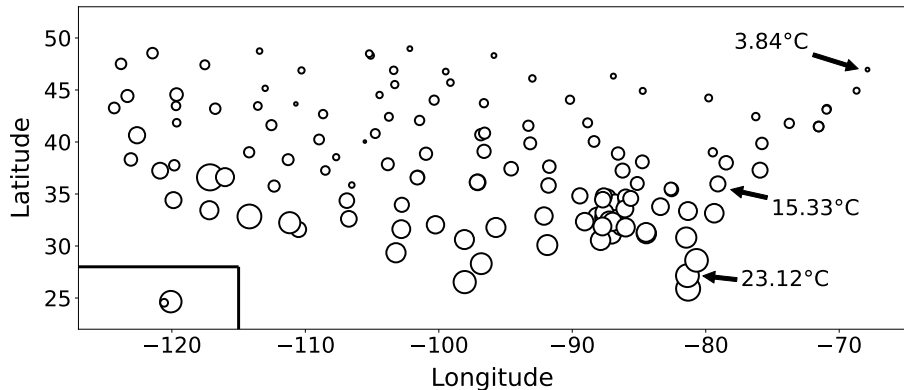
The sample mean of $X := \{x_1, x_2, \dots, x_n\}$ is

$$m(X) := \frac{\sum_{i=1}^n x_i}{n}$$

Temperature dataset

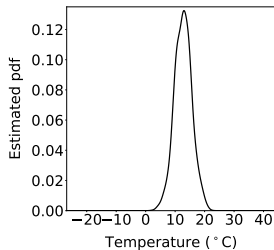
Hourly temperatures at 134 weather stations in the US

○ Weather-station locations (radius proportional to mean temperature)

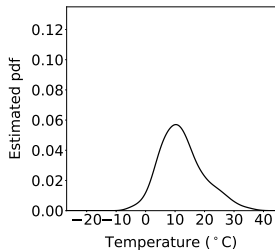


Same mean

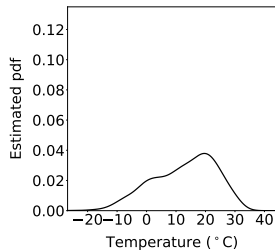
$$m(X) = 12.7^{\circ}\text{C}$$



$$m(X) = 12.3^{\circ}\text{C}$$



$$m(X) = 12.7^{\circ}\text{C}$$



Challenge

Quantifying *magnitude* of deviation from the mean

Magnitude of a random variable?

Magnitude of real number a : $|a| = \sqrt{a^2}$

Euclidean length of vector x : $\|x\|_2 = \sqrt{\sum_{i=1}^d x[i]^2}$

Magnitude/energy of random variable \tilde{a} ?

Mean square or second moment $\mathbb{E} [\tilde{a}^2]$

Mean squared error

The mean squared error (MSE) between an estimate \tilde{e} and a random variable \tilde{a} is

$$\text{E} [(\tilde{a} - \tilde{e})^2]$$

Minimum MSE constant estimate

Best **constant** estimate of \tilde{a} ?

$$\arg \min_{c \in \mathbb{R}} \mathbb{E} [(c - \tilde{a})^2] = \mathbb{E}[\tilde{a}]$$

$$\text{MSE}(c) := \mathbb{E} [(c - \tilde{a})^2] = c^2 - 2c\mathbb{E}[\tilde{a}] + \mathbb{E}[\tilde{a}^2]$$

$$\text{MSE}'(c) = 2(c - \mathbb{E}[\tilde{a}])$$

$$\text{MSE}''(c) = 2$$

The **mean** $\mathbb{E}[\tilde{a}]$

Variance

Mean squared distance of a random variable to its mean

$$\begin{aligned}\text{Var} [\tilde{a}] &:= \text{E} \left[(\tilde{a} - \text{E} [\tilde{a}])^2 \right] \\ &= \text{E} \left[\tilde{a}^2 - 2\tilde{a}\text{E} [\tilde{a}] + \text{E} [\tilde{a}]^2 \right] \\ &= \text{E} [\tilde{a}^2] - 2\text{E} [\tilde{a}] \text{E} [\tilde{a}] + \text{E} [\tilde{a}]^2 \\ &= \text{E} [\tilde{a}^2] - \text{E}(\tilde{a})^2\end{aligned}$$

Standard deviation

The standard deviation $\sigma_{\tilde{a}}$ of \tilde{a} is

$$\sigma_{\tilde{a}} := \sqrt{\text{Var}[\tilde{a}]}$$

Sample variance

Dataset: x_1, x_2, \dots, x_n

The sample variance is the average squared deviation from the sample mean

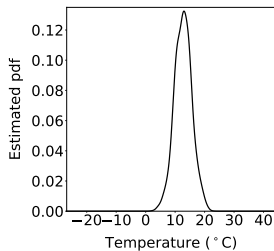
$$v(X) := \frac{\sum_{i=1}^n (x_i - m(X))^2}{n - 1}$$

The sample standard deviation σ_X is the square root of the sample variance

Same mean

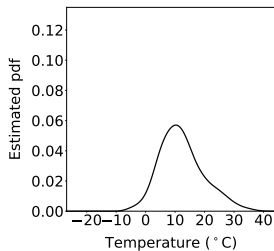
$$m(X) = 12.7^{\circ}\text{C}$$

$$\sqrt{v(X)} = 2.9^{\circ}\text{C}$$



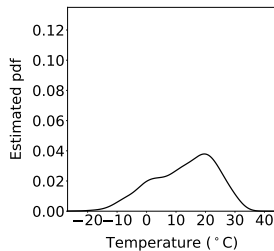
$$m(X) = 12.3^{\circ}\text{C}$$

$$\sqrt{v(X)} = 7.5^{\circ}\text{C}$$

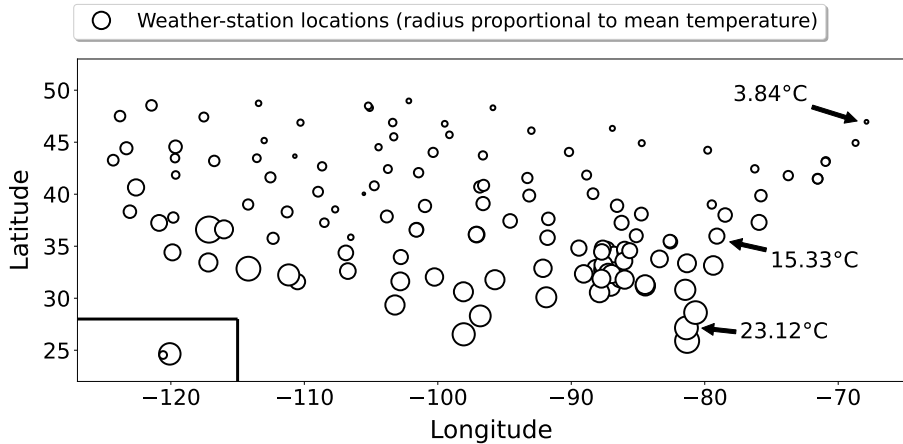


$$m(X) = 12.7^{\circ}\text{C}$$

$$\sqrt{v(X)} = 10.6^{\circ}\text{C}$$

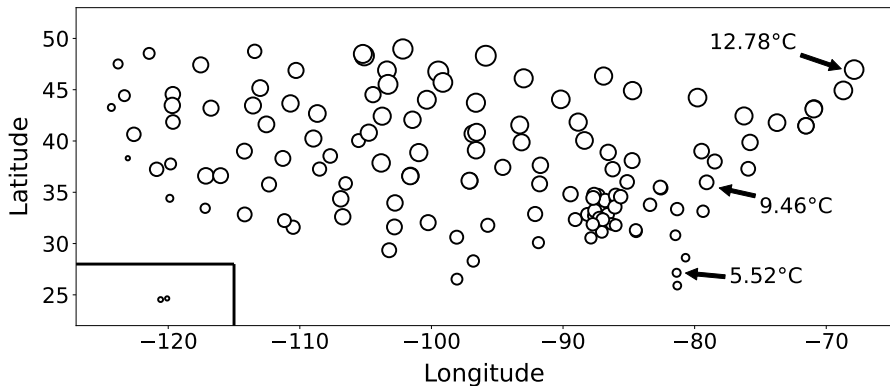


Means



Standard deviations

○ Weather-station locations (radius proportional to standard deviation of temperature)



Scaled, shifted variable

We scale a random variable by c_1 and shift it by c_2

Variance?

$$\begin{aligned}\text{Var}[c_1 \tilde{a} + c_2] &= \text{E} \left[(c_1 \tilde{a} + c_2 - \text{E}[c_1 \tilde{a} + c_2])^2 \right] \\ &= \text{E} \left[(c_1 \tilde{a} + c_2 - c_1 \text{E}[\tilde{a}] - c_2)^2 \right] \\ &= c_1^2 \text{E} \left[(\tilde{a} - \text{E}[\tilde{a}])^2 \right] \\ &= c_1^2 \text{Var}[\tilde{a}]\end{aligned}$$

Bernoulli random variable

$$\mathbb{E}[\tilde{a}] = \theta$$

$$\begin{aligned}\mathbb{E}[\tilde{a}^2] &= 0 \cdot p_{\tilde{a}}(0) + 1 \cdot p_{\tilde{a}}(1) \\ &= \theta\end{aligned}$$

$$\begin{aligned}\text{Var}[\tilde{a}] &= \mathbb{E}[\tilde{a}^2] - \mathbb{E}[\tilde{a}]^2 \\ &= \theta(1 - \theta)\end{aligned}$$

Geometric random variable

$$\sum_{k=1}^{\infty} k \alpha^k = \frac{\alpha}{(1-\alpha)^2}$$

$$\sum_{k=1}^{\infty} k^2 \alpha^{k-1} = \frac{1+\alpha}{(1-\alpha)^3}$$

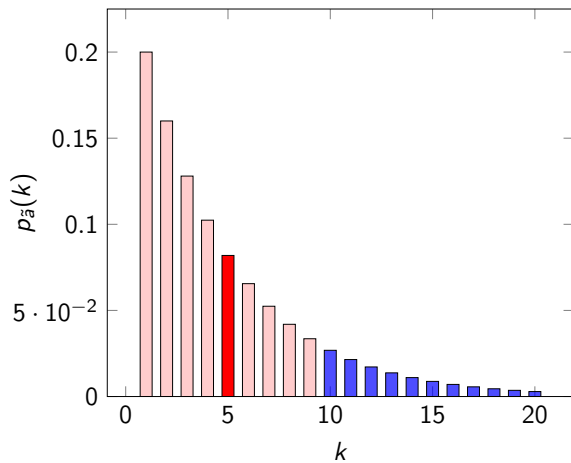
$$\begin{aligned} \mathbb{E} [\tilde{a}^2] &= \sum_{k=1}^{\infty} k^2 p_{\tilde{a}}(k) \\ &= \sum_{k=1}^{\infty} k^2 \theta (1-\theta)^{k-1} \\ &= \frac{2-\theta}{\theta^2} \end{aligned}$$

Geometric random variable

$$\begin{aligned} \mathbb{E}[\tilde{a}] &= \frac{1}{\theta} \\ \mathbb{E}[\tilde{a}^2] &= \frac{2 - \theta}{\theta^2} \end{aligned}$$

$$\begin{aligned} \text{Var}[\tilde{a}] &= \mathbb{E}[\tilde{a}^2] - \mathbb{E}[\tilde{a}]^2 \\ &= \frac{1 - \theta}{\theta^2} \end{aligned}$$

Geometric, $\theta := 0.2$



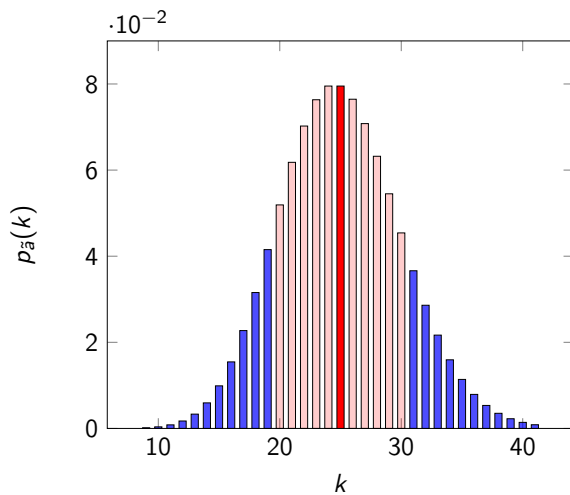
Poisson random variable

$$\sum_{k=0}^{\infty} \frac{\lambda^k}{k!} = e^{\lambda}$$

$$\begin{aligned} \mathbb{E} [\tilde{a}^2] &= \sum_{k=1}^{\infty} k^2 p_{\tilde{a}}(k) = \sum_{k=1}^{\infty} \frac{k^2 \lambda^k e^{-\lambda}}{k!} \\ &= \sum_{k=1}^{\infty} \frac{k \lambda^k e^{-\lambda}}{(k-1)!} \\ &= e^{-\lambda} \left(\sum_{k=1}^{\infty} \frac{(k-1) \lambda^k}{(k-1)!} + \frac{\lambda^k}{(k-1)!} \right) \\ &= e^{-\lambda} \left(\sum_{m=0}^{\infty} \frac{\lambda^{m+2}}{m!} + \sum_{m=0}^{\infty} \frac{\lambda^{m+1}}{m!} \right) \\ &= \lambda^2 + \lambda \end{aligned}$$

$$\text{Var} [\tilde{a}] = \mathbb{E} [\tilde{a}^2] - \mathbb{E} [\tilde{a}]^2 = \lambda$$

Poisson, $\lambda := 25$

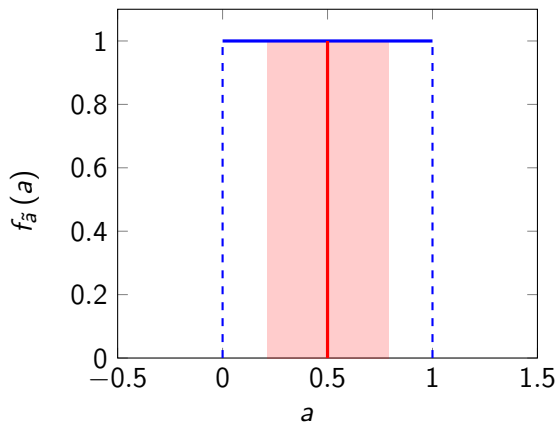


Uniform random variable

$$\begin{aligned} \mathbb{E} [\tilde{u}] &= \frac{a + b}{2} \\ \mathbb{E} [\tilde{u}^2] &= \int_{u=-\infty}^{\infty} u^2 f_{\tilde{u}}(u) \, du \\ &= \int_{u=a}^b \frac{u^2}{b-a} \, du \\ &= \frac{b^3 - a^3}{3(b-a)} \\ &= \frac{a^2 + ab + b^2}{3} \end{aligned}$$

$$\begin{aligned} \text{Var} [\tilde{u}] &= \mathbb{E} [\tilde{u}^2] - \mathbb{E} [\tilde{u}]^2 \\ &= \frac{(b-a)^2}{12} \end{aligned}$$

Uniform in $[0, 1]$

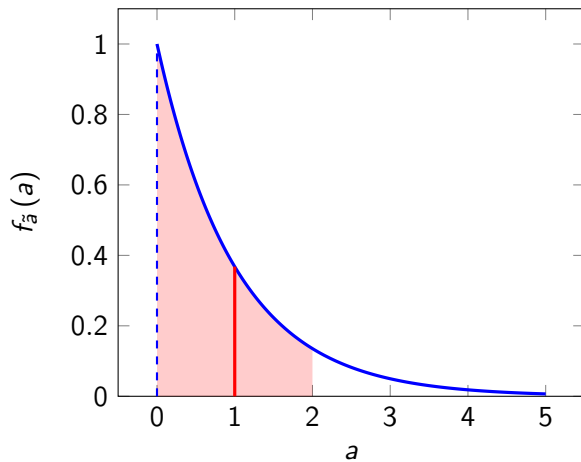


Exponential random variable

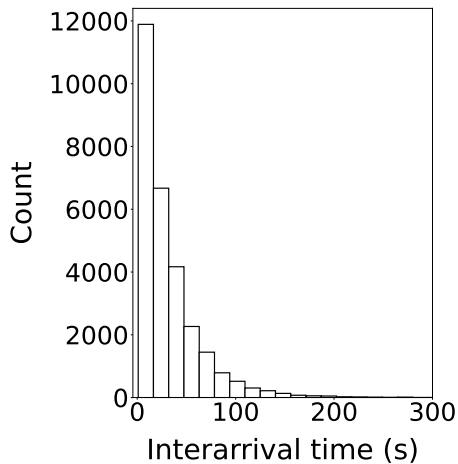
$$\begin{aligned} \mathrm{E} [\tilde{a}] &= \frac{1}{\lambda} \\ \mathrm{E} [\tilde{a}^2] &= \int_{a=-\infty}^{\infty} a^2 f_{\tilde{a}}(a) \, da \\ &= \int_{a=0}^{\infty} a^2 \lambda e^{-\lambda a} \, da \\ &= a^2 e^{-\lambda a} \Big|_0^{\infty} + 2 \frac{1}{\lambda} \int_0^{\infty} a \lambda e^{-\lambda a} \, da \\ &= \frac{2}{\lambda^2} \end{aligned}$$

$$\mathrm{Var} [\tilde{a}] = \mathrm{E} [\tilde{a}^2] - \mathrm{E} [\tilde{a}]^2 = \frac{1}{\lambda^2}$$

Exponential, $\lambda := 1$



Call center data



Sample mean = 30.8

Sample standard deviation = 33.6

Gaussian random variable

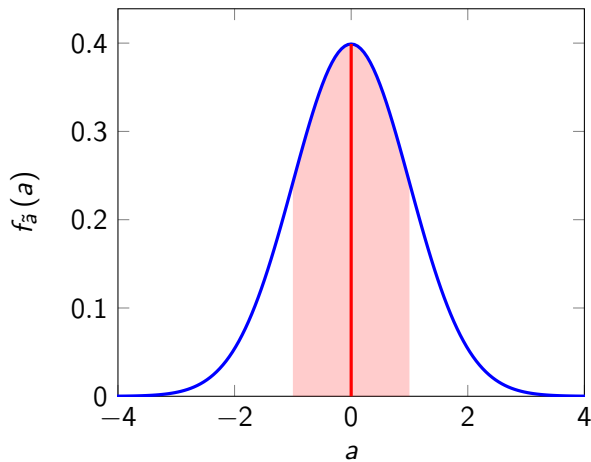
Change of variables $t = (a - \mu) / \sigma$

$$\mathbb{E} [\tilde{a}] = \mu$$

$$\begin{aligned}\mathbb{E} [\tilde{a}^2] &= \int_{a=-\infty}^{\infty} a^2 f_{\tilde{a}}(a) \, da \\&= \int_{a=-\infty}^{\infty} \frac{a^2}{\sqrt{2\pi}\sigma} e^{-\frac{(a-\mu)^2}{2\sigma^2}} \, da \\&= \frac{\sigma^2}{\sqrt{2\pi}} \int_{t=-\infty}^{\infty} t^2 e^{-\frac{t^2}{2}} \, dt + \frac{2\mu\sigma}{\sqrt{2\pi}} \int_{t=-\infty}^{\infty} t e^{-\frac{t^2}{2}} \, dt \\&\quad + \frac{\mu^2}{\sqrt{2\pi}} \int_{t=-\infty}^{\infty} e^{-\frac{t^2}{2}} \, dt \\&= \frac{\sigma^2}{\sqrt{2\pi}} \left(t^2 e^{-\frac{t^2}{2}} \Big|_{-\infty}^{\infty} + \int_{t=-\infty}^{\infty} e^{-\frac{t^2}{2}} \, dt \right) + \mu^2 \\&= \sigma^2 + \mu^2\end{aligned}$$

$$\text{Var} [\tilde{a}] = \mathbb{E} [\tilde{a}^2] - \mathbb{E} [\tilde{a}]^2 = \sigma^2$$

Gaussian $\mu := 0$, $\sigma^2 := 1$



What have we learned

Definition of mean square, variance and standard deviation

Variance of popular parametric models