

Variance

Lemma 1: Suppose X_1, \dots, X_n be a sample (i.i.d., SRS, or other). Then

$$\sum_i (X_i - \bar{X})^2 = \dots = \sum_i X_i^2 - n\bar{X}^2.$$

Proof: Expand the expression. □

Let $\mathbf{X} = \langle X_1, \dots, X_n \rangle$, $\bar{\mathbf{X}} = \bar{X}\mathbf{1} = \bar{X} \langle 1, \dots, 1 \rangle = \langle \bar{X}, \dots, \bar{X} \rangle$, and $\mathbf{V} = \mathbf{X} - \bar{\mathbf{X}}$. Note that

$$|\mathbf{X}|^2 = \sum_i X_i^2, \quad |\bar{\mathbf{X}}|^2 = n\bar{X}^2, \quad \text{and} \quad |\mathbf{V}|^2 = \sum_i (X_i - \bar{X})^2.$$

Lemma 2: Suppose X_1, \dots, X_n is an i.i.d. random sample from a population with mean μ , variance σ^2 . Then the sample variance, defined as

$$S^2 := \frac{\sum_i (X_i - \bar{X})^2}{n-1}, \quad = \frac{|\vec{V}|^2}{n-1}$$

is an unbiased estimator for the population variance σ^2 . That is,

$$E\left(\frac{\sum_i (X_i - \bar{X})^2}{n-1}\right) = \sigma^2.$$

Proof: First, note in general that, since $\text{Var}(W) = E(W^2) - [E(W)]^2$, we have that

$$E(X_i^2) = \mu^2 + \sigma^2, \quad \text{for each } i, \text{ and} \quad E(\bar{X}^2) = \mu^2 + \frac{\sigma^2}{n}.$$

Combining this with the result of the previous lemma, we show

$$E\left(\sum_i (X_i - \bar{X})^2\right) = E\left(\sum_i X_i^2 - n\bar{X}^2\right) = \sum_i E(X_i^2) - nE(\bar{X}^2) = \dots = (n-1)\sigma^2,$$

which leads to the claim. □

Notes:

- For the seemingly "more natural" estimator

$$\hat{\sigma}^2 = \frac{\sum_i X_i - \bar{X})^2}{n}, \quad \text{we have} \quad E(\hat{\sigma}^2) = \frac{n-1}{n}\sigma^2.$$

That is, it is *biased*.

$$\vec{X} = \langle x_1, x_2, \dots, x_n \rangle, \quad \bar{X} = \frac{1}{n} \sum x_i$$

- Viewing the results of the lemmas above in terms of vector lengths, we have

$$(n-1)S^2 = |\mathbf{V}|^2 = |\mathbf{X}|^2 - |\bar{\mathbf{X}}|^2,$$

or

$$|\mathbf{V}|^2 + |\bar{\mathbf{X}}|^2 = |\mathbf{X}|^2 = |\mathbf{V} + \bar{\mathbf{X}}|^2.$$

Since the Pythagorean Theorem is an "if and only if" result, it follows that $\mathbf{V} \perp \bar{\mathbf{X}}$.

- We may choose an orthonormal basis $\{\mathbf{u}_1, \mathbf{u}_2, \dots, \mathbf{u}_n\}$ for \mathbb{R}^n with \mathbf{u}_1 parallel to $\mathbf{1}$. Then

$$\mathbf{X} = \bar{\mathbf{X}} + \mathbf{V}$$

$$\mathcal{U} = \text{span}\{\bar{\mathbf{u}}_2, \dots, \bar{\mathbf{u}}_n\}$$

where

- $\bar{\mathbf{X}} = \text{proj}(\mathbf{X} \rightarrow \mathbf{u}_1)$
- If we write $\mathcal{U} = \text{span}(\mathbf{u}_2, \dots, \mathbf{u}_n)$, then

$$\mathbf{V} = \text{proj}(\mathbf{X} \rightarrow \mathcal{U}) = \sum_{i=2}^n \text{proj}(\mathbf{X} \rightarrow \mathbf{u}_i) = \sum_{i=2}^n (\mathbf{X} \cdot \mathbf{u}_i) \mathbf{u}_i,$$

the second equation a result of orthogonality (see Friday's notes). And, after repeated application of the Pythagorean theorem,

$$(n-1)S^2 = |\mathbf{V}|^2 = \sum_{i=2}^n (\mathbf{X} \cdot \mathbf{u}_i)^2,$$

- If $\langle X_1, \dots, X_n \rangle \stackrel{\text{i.i.d.}}{\sim} \text{Norm}(\mu, \sigma)$, then by defining $W_i := \mathbf{X} \cdot \mathbf{u}_i$ for $i = 2, \dots, n$, we get the $W_i \stackrel{\text{i.i.d.}}{\sim} \text{Norm}(0, \sigma)$, and

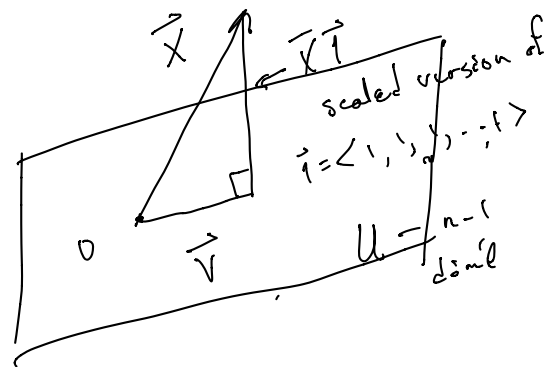
$$(n-1)S^2 = |\mathbf{V}|^2 = \sum_{i=2}^n (\mathbf{X} \cdot \mathbf{u}_i)^2 = W_2^2 + W_3^2 + \dots + W_n^2.$$

In light of all this, we begin using a new distribution.

Definition 1: Let $\mathbf{Z} = \langle Z_1, \dots, Z_n \rangle \stackrel{\text{i.i.d.}}{\sim} \text{Norm}(0, 1)$. The sum

$$Z_1^2 + Z_2^2 + \dots + Z_n^2$$

has the **chi-squared distribution** with n **degrees of freedom**, abbreviated as **Chisq**(n).



Using the cdf method, for $X = Z_1^2 + \dots + Z_n^2$ as above, one can show that $X \sim \text{Chisq}(n)$, then $X \sim \text{Gamma}(\alpha = n/2, \lambda = 1/2)$

Case $(n=1)$: $X = Z_1^2$ w/ $Z \sim \text{Norm}(0, 1)$

$$\frac{F_X(x)}{F_X(x)} = P_r(X \leq x) = P_r(Z_1^2 \leq x) = P_r(-\sqrt{x} \leq Z_1 \leq \sqrt{x})$$

$$= \Phi(\sqrt{x}) - \Phi(-\sqrt{x})$$

So $f_X(x) = \frac{d}{dx} [\Phi(\sqrt{x}) - \Phi(-\sqrt{x})] = \frac{1}{\sqrt{2\pi}} e^{-(\sqrt{x})^2/2} \cdot \frac{1}{2\sqrt{x}} + \frac{1}{\sqrt{2\pi}} e^{-(-\sqrt{x})^2/2} \cdot \frac{1}{2\sqrt{x}}$

$$= \frac{1}{\sqrt{2\pi}} \cdot \frac{1}{2\sqrt{x}} e^{-x/2} (1 + 1) = \left(\frac{1}{\sqrt{2\pi}} \cdot \frac{1}{\sqrt{x}} e^{-x/2} \right)$$

= pdf for another dist., namely $\text{Gamma}(\alpha = 1/2, \lambda = 1/2)$

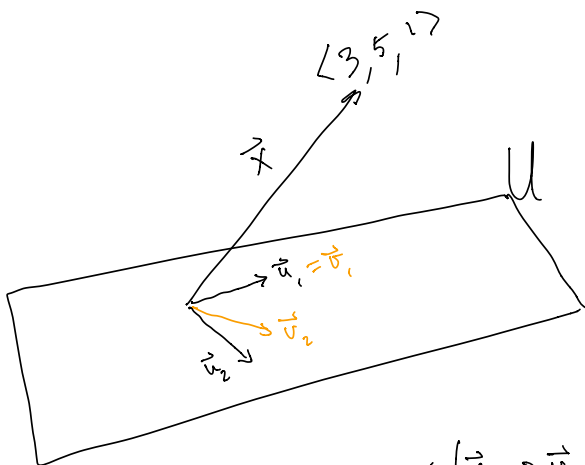
Generally: $\underbrace{Z_1^2 + \dots + Z_n^2}_{\text{independent}} \sim \text{Gamma}(n/2, 1/2)$

So, $\text{Chisq}(n)$ is a special case of a Gamma distribution. We have introduced it because of this result:

Lemma 3 (Lemma 4.6.6, p. 267): Let $\mathbf{X} = \langle X_1, \dots, X_n \rangle \stackrel{\text{i.i.d.}}{\sim} \text{Norm}(\mu, \sigma)$. Then

(i) $\frac{(n-1)S^2}{\sigma^2} = \frac{\sum_{i=1}^n (X_i - \bar{X})^2}{\sigma^2} \sim \text{Chisq}(n-1)$, and

(ii) \bar{X} and S^2 are independent random variables.



$$\begin{aligned}
 & \text{proj}(\vec{x} \rightarrow \vec{u}_1) + \text{proj}(\vec{x} \rightarrow \vec{u}_2) \\
 &= \langle 4, 4, 0 \rangle + \langle -1, 1, 0 \rangle \\
 &= \langle 3, 5, 0 \rangle \\
 &\stackrel{?}{=} \text{proj}(\vec{x} \rightarrow U)
 \end{aligned}$$

$$\begin{aligned}
 & \text{proj}(\vec{x} \rightarrow \vec{v}_1) + \text{proj}(\vec{x} \rightarrow \vec{v}_2) \\
 &= \langle 4, 4, 0 \rangle + \langle 3, 0, 0 \rangle = \langle 7, 4, 0 \rangle \\
 &\stackrel{?}{=} \text{proj}(\vec{x} \rightarrow U)
 \end{aligned}$$

Conclusion(?): When you have an orthogonal basis for a subspace U , then

$$\text{proj}(\vec{x} \rightarrow U) = \text{sum of projections} \left(\vec{u} \rightarrow \begin{matrix} \text{basis} \\ \text{vectors} \end{matrix} \right).$$