

Length n Probability Generation

Matthew Seguin

Summary

The goals of this document are:

- Explain what an probability vector of length n is and what properties such a vector has.
- Illustrate how to generate a random probability vector of length n .
- Analyze how to bias such a probability vector.
- Examine the distribution of the norms of such vectors.

Importing Libraries

```
library(tidyverse)
library(latex2exp)
```

Generator

```
prob_generator <- function(N, m = 1000, shape = NA){  
  shape <- ifelse(is.na(shape), rep(1, N), shape)  
  gams <- matrix(rgamma(n = N*m, shape = shape), nrow = N)  
  probs <- t(t(gams)/colSums(gams))  
  return(probs)  
}
```

This just uses the Dirichlet distribution for generating the probabilities.

The intuition is that we generate $X_j \sim \text{Gamma}(\alpha_j, \lambda)$ for $j \in \{1, \dots, n\}$. Then we consider the distribution of

$$p = (p_1, \dots, p_n) = \left(\frac{X_1}{S}, \dots, \frac{X_n}{S} \right) \text{ where } S = \sum_{j=1}^n X_j.$$

Note that this is an $n - 1$ dimensional vector because we know $\sum_{j=1}^n p_j = \sum_{j=1}^n \frac{X_j}{S} = \frac{1}{S} \sum_{j=1}^n X_j = \frac{S}{S} = 1$. Quick note:

once we know $n - 1$ of the p_j we know the last one.

Furthermore we also know $X_j \geq 0$ for each $j \in \{1, \dots, n\}$ (since X_j follows a gamma distribution) so

$S = \sum_{j=1}^n X_j > X_j \geq 0$ for each $j \in \{1, \dots, n\}$ as well. Which then tells us that $p_j = \frac{X_j}{S} \geq 0$ and $p_j = \frac{X_j}{S} \leq 1$.

So $p = (p_1, \dots, p_n)$ is indeed a probability vector of length n .

The support of the Dirichlet distribution is the standard $k - 1$ simplex. Which is given by

$$\Delta_{k-1} = \left\{ (p_1, \dots, p_n) \in \mathbb{R}^n : \forall_{j \in \{1, \dots, n\}} 0 \leq p_j \leq 1, \sum_{j=1}^n p_j = 1 \right\}.$$

Furthermore the density can be given by:

$$f_{p_1, \dots, p_n}(p_1, \dots, p_n) = \frac{1}{B(\alpha)} \prod_{j=1}^n x_j^{\alpha_j - 1} \text{ for } (p_1, \dots, p_n) \in \Delta_{k-1}$$

$$\text{Where } B(\alpha) = \frac{\Gamma(\alpha_1) \dots \Gamma(\alpha_n)}{\Gamma(\alpha_1 + \dots + \alpha_n)}.$$

The special case where $\alpha_j = 1$ for each $j \in \{1, \dots, n\}$ is essentially a uniform distribution over the standard $k - 1$ simplex

as the density is given by:

$$f_{p_1, \dots, p_n}(p_1, \dots, p_n) = \frac{1}{B(\alpha)} \prod_{j=1}^n x_j^{\alpha_j - 1} = \frac{1}{B(\alpha)} \prod_{j=1}^n x_j^{1-1} = \frac{1}{B(\alpha)} = \frac{\Gamma(\alpha_1 + \dots + \alpha_n)}{\Gamma(\alpha_1) \dots \Gamma(\alpha_n)}$$

$$= \frac{\Gamma(n)}{(\Gamma(1))^n} = \frac{(n-1)!}{(0!)^n} = (n-1)! \text{ for } (p_1, \dots, p_n) \in \Delta_{k-1}$$

It is uniform since clearly the density does not depend on $p = (p_1, \dots, p_n)$.

The above result uses the fact that $\Gamma(k) = (k-1)!$ for $k \in \{1, 2, \dots\}$ which I will prove below:

$$\Gamma(k) = \int_0^\infty t^{k-1} e^{-t} dt$$

This clearly converges for $k-1 \geq 0$ i.e. for $k \geq 1$.

Now let $u(t) = t^{k-1}$ and $\frac{dv}{dt} = e^{-t}$ so that $\frac{du}{dt} = (k-1)t^{k-2}$ and $v(t) = -e^{-t}$. Then:

$$\begin{aligned} \Gamma(k) &= \int_0^\infty t^{k-1} e^{-t} dt = u(t)v(t) \Big|_0^\infty - \int_0^\infty v(t) \frac{du}{dt} dt = -t^{k-1} e^{-t} \Big|_0^\infty + \int_0^\infty (k-1)t^{k-2} e^{-t} dt \\ &= \lim_{t \rightarrow \infty} -\frac{t^{k-1}}{e^t} + 0 + (k-1) \int_0^\infty t^{(k-1)-1} e^{-t} dt = (k-1)\Gamma(k-1) \end{aligned}$$

Which again converges for $k-1 \geq 1$ or equivalently $k \geq 2$.

Therefore if $k \geq 2$ then $\Gamma(k) = (k-1)\Gamma(k-1)$ which seems like a factorial, now we need our base case:

$$\Gamma(1) = \int_0^\infty t^{1-1} e^{-t} dt = \int_0^\infty e^{-t} dt = \int_{-\infty}^\infty f_T(t) dt = 1 = 0! \quad \text{where } T \sim \text{Exponential}(1)$$

Therefore if $k \geq 1$ we know for

$$\Gamma(k) = (k-1)\Gamma(k-1) = \dots = (k-1)\dots(2)\Gamma(2) = (k-1)\dots(2)(1)\Gamma(1) = (k-1)\dots(2)(1) = (k-1)! \quad \square$$

Testing that this does indeed produce length n probability vectors.

There is always the issue of machine precision for floating point numbers so the sums can not be exactly one but we can make a simple tolerance based on machine epsilon.

```
tol <- 4.75*.Machine$double.eps
tol
```

```
## [1] 1.054712e-15
```

```
for (i in 2:1000){
  test <- prob_generator(N = i)
  stopifnot(
    min(test) >= 0,
    max(abs(colSums(test) - 1)) <= tol
  )
}
```

As there are no errors we can see that all of these have non-negative probabilities and that they sum to 1 (within our tolerance based on machine imprecision) and hence define length n probability vectors.

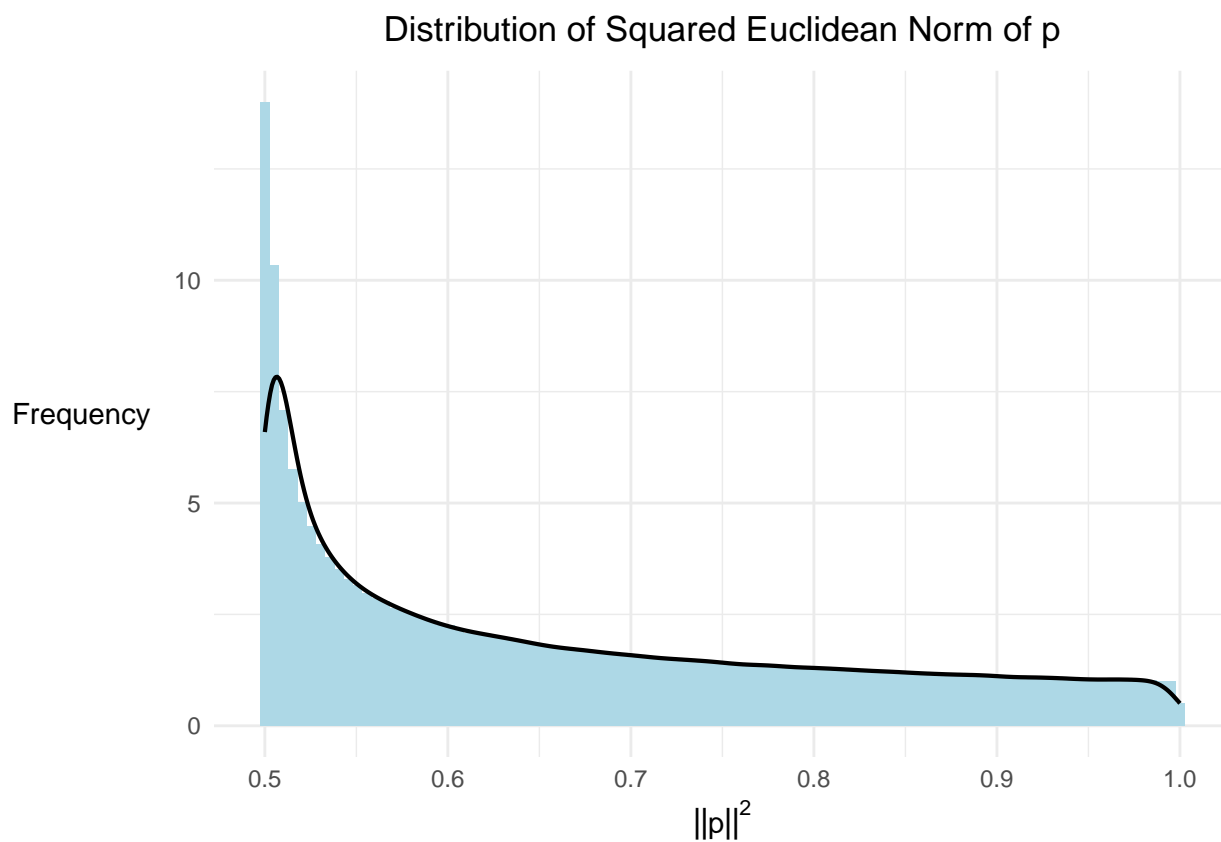
Examining Distribution of the Square of the Euclidean Norm

Below are some graphs to analyze the distribution of $\|p\|^2$.

- First $n = 2$:

```
probs <- prob_generator(N = 2, m = 1000000)
sq_norms <- data.frame(sq_norm = colSums(probs^2))

sq_norms %>%
  ggplot(aes(x = sq_norm)) +
  geom_histogram(aes(y = after_stat(density)),
    fill = "lightblue",
    bins = 100) +
  geom_density(col = "black",
    linewidth = 0.75) +
  labs(x = TeX("$\|p\|^2$"),
    y = "Frequency",
    title = TeX("Distribution of Squared Euclidean Norm of $p$")) +
  theme_minimal() +
  theme(plot.title = element_text(hjust = 0.5),
    axis.title = element_text(color = "black"),
    axis.title.y = element_text(angle = 0, vjust = 0.5))
```

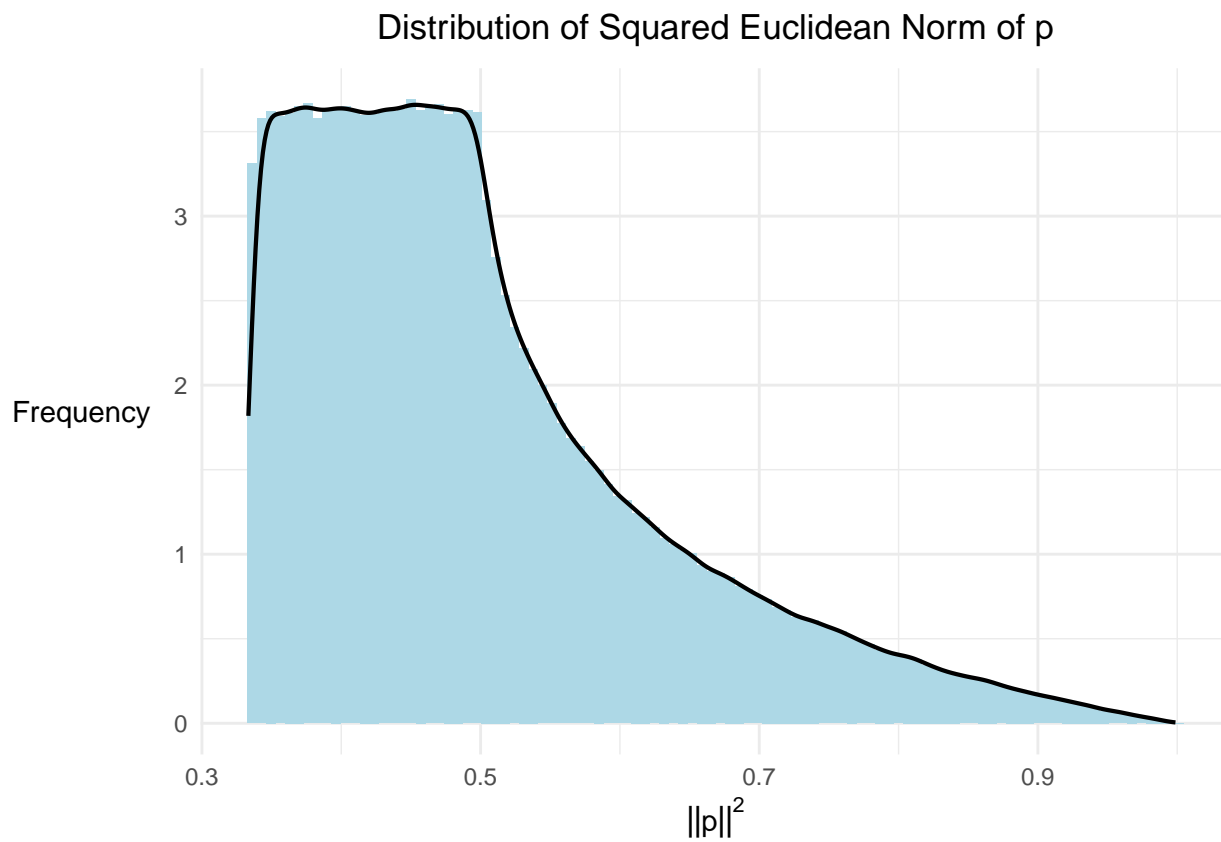


Comment.

- Now $n = 3$:

```
probs <- prob_generator(N = 3, m = 1000000)
sq_norms <- data.frame(sq_norm = colSums(probs^2))

sq_norms %>%
  ggplot(aes(x = sq_norm)) +
  geom_histogram(aes(y = after_stat(density)),
    fill = "lightblue",
    bins = 100) +
  geom_density(col = "black",
    linewidth = 0.75) +
  labs(x = TeX("$||p||^2$"),
    y = "Frequency ",
    title = TeX("Distribution of Squared Euclidean Norm of $p$")) +
  theme_minimal() +
  theme(plot.title = element_text(hjust = 0.5),
    axis.title = element_text(color = "black"),
    axis.title.y = element_text(angle = 0, vjust = 0.5))
```

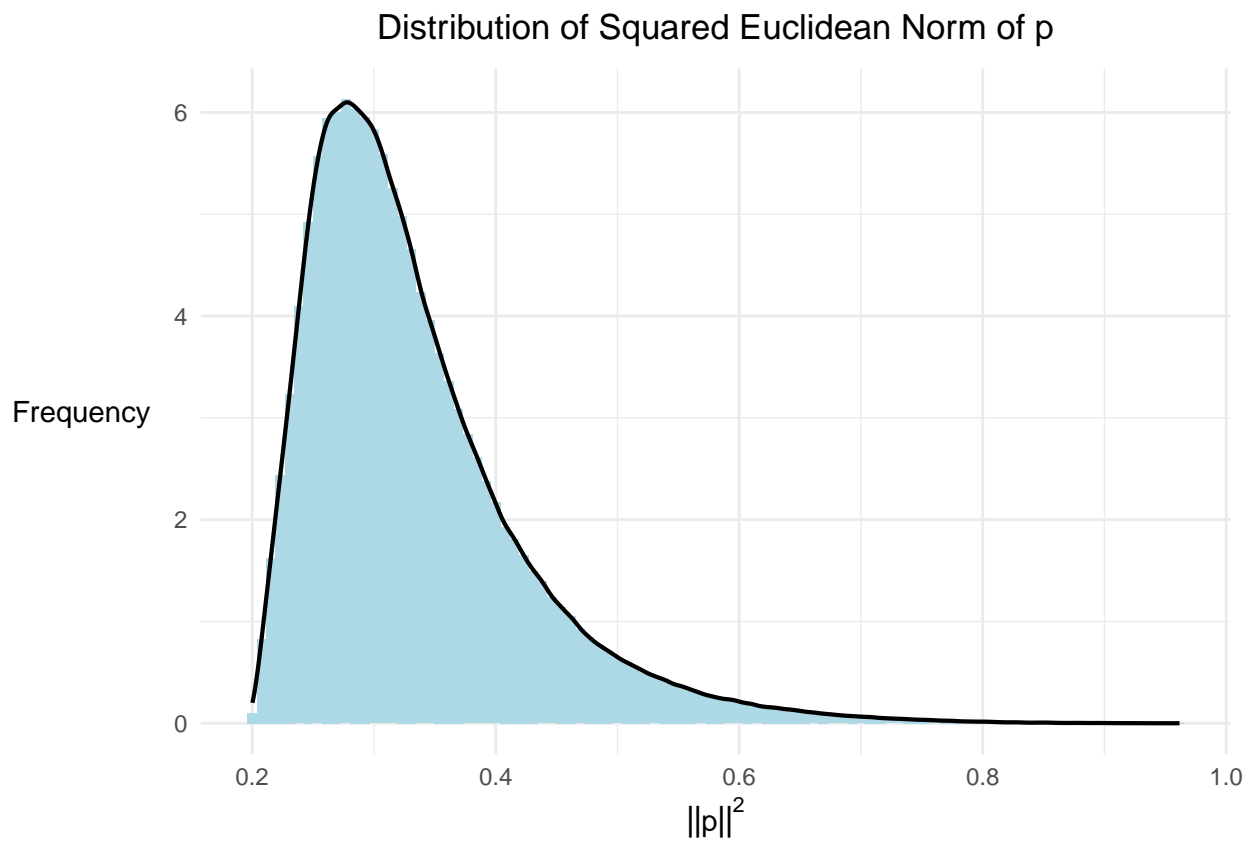


Comment.

- Now $n = 5$:

```
probs <- prob_generator(N = 5, m = 1000000)
sq_norms <- data.frame(sq_norm = colSums(probs^2))

sq_norms %>%
  ggplot(aes(x = sq_norm)) +
  geom_histogram(aes(y = after_stat(density)),
    fill = "lightblue",
    bins = 100) +
  geom_density(col = "black",
    linewidth = 0.75) +
  labs(x = TeX("$||p||^2$"),
    y = "Frequency ",
    title = TeX("Distribution of Squared Euclidean Norm of $p$")) +
  theme_minimal() +
  theme(plot.title = element_text(hjust = 0.5),
    axis.title = element_text(color = "black"),
    axis.title.y = element_text(angle = 0, vjust = 0.5))
```

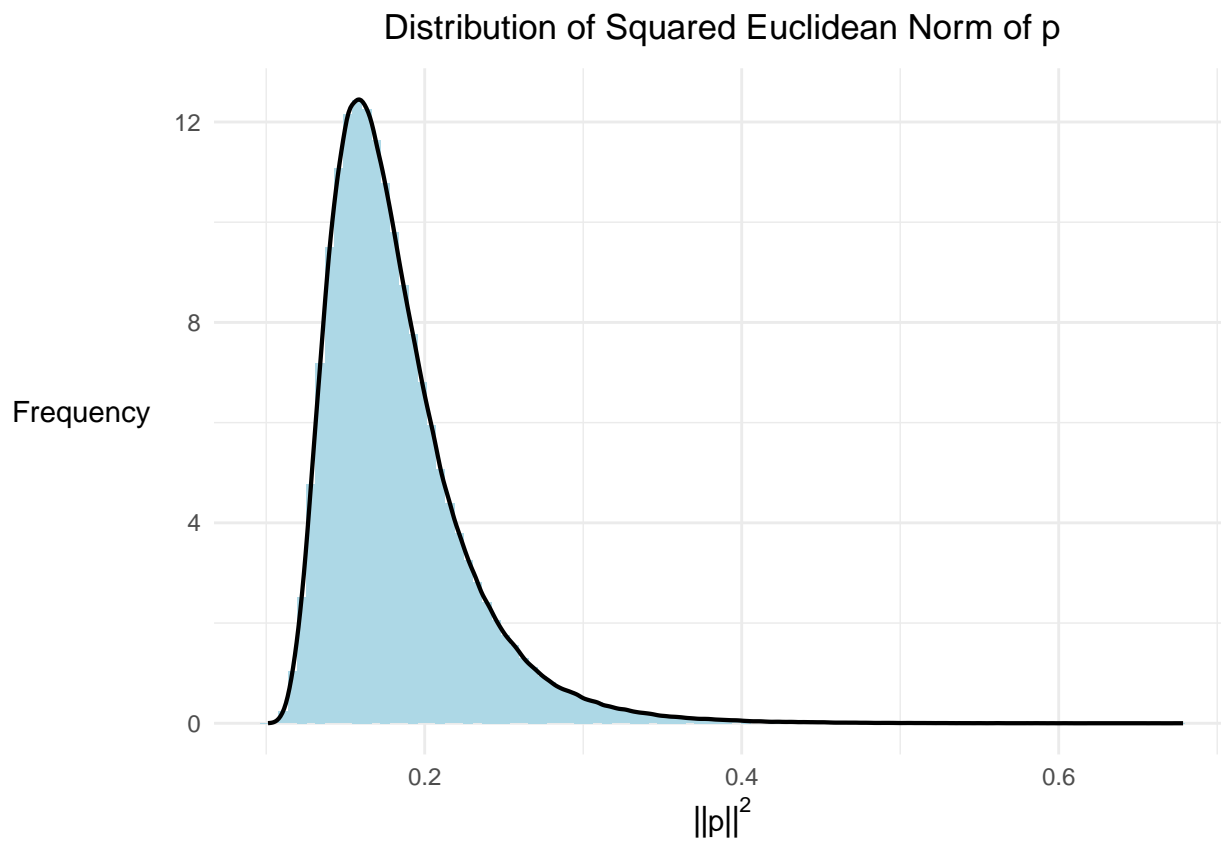


Comment.

- Now $n = 10$:

```
probs <- prob_generator(N = 10, m = 1000000)
sq_norms <- data.frame(sq_norm = colSums(probs^2))

sq_norms %>%
  ggplot(aes(x = sq_norm)) +
  geom_histogram(aes(y = after_stat(density)),
    fill = "lightblue",
    bins = 100) +
  geom_density(col = "black",
    linewidth = 0.75) +
  labs(x = TeX("$||p||^2$"),
    y = "Frequency ",
    title = TeX("Distribution of Squared Euclidean Norm of $p$")) +
  theme_minimal() +
  theme(plot.title = element_text(hjust = 0.5),
    axis.title = element_text(color = "black"),
    axis.title.y = element_text(angle = 0, vjust = 0.5))
```

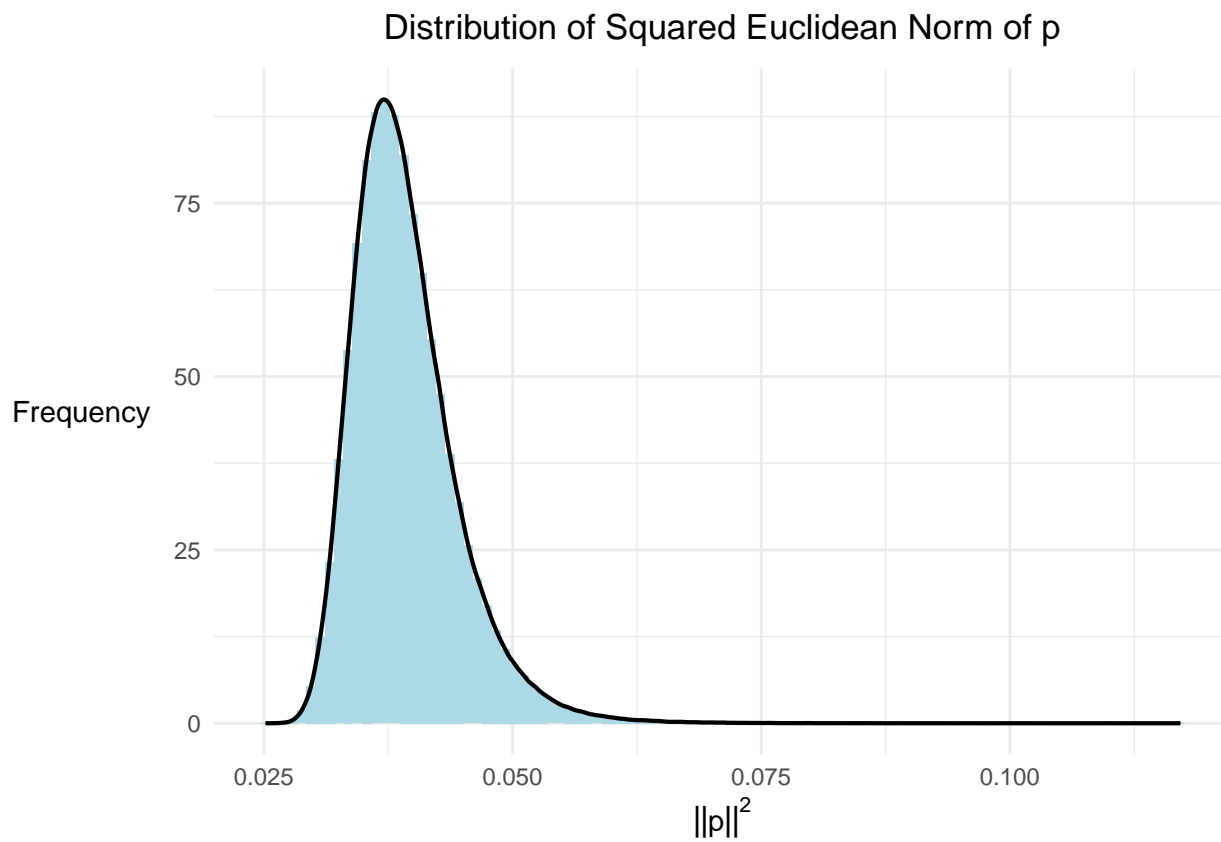


Comment.

- Now $n = 50$:

```
probs <- prob_generator(N = 50, m = 1000000)
sq_norms <- data.frame(sq_norm = colSums(probs^2))

sq_norms %>%
  ggplot(aes(x = sq_norm)) +
  geom_histogram(aes(y = after_stat(density)),
    fill = "lightblue",
    bins = 100) +
  geom_density(col = "black",
    linewidth = 0.75) +
  labs(x = TeX("$||p||^2$"),
    y = "Frequency ",
    title = TeX("Distribution of Squared Euclidean Norm of $p$")) +
  theme_minimal() +
  theme(plot.title = element_text(hjust = 0.5),
    axis.title = element_text(color = "black"),
    axis.title.y = element_text(angle = 0, vjust = 0.5))
```

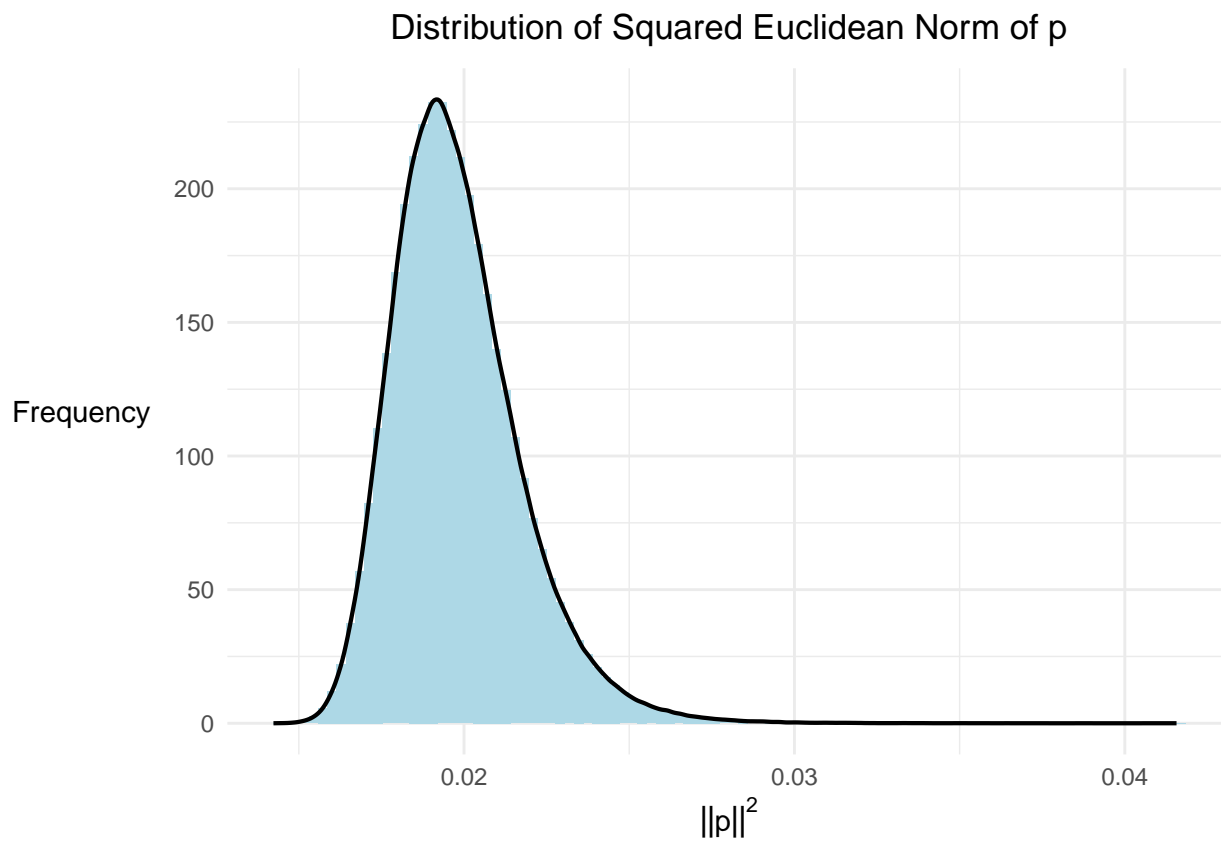


Comment.

- Now $n = 100$:

```
probs <- prob_generator(N = 100, m = 1000000)
sq_norms <- data.frame(sq_norm = colSums(probs^2))

sq_norms %>%
  ggplot(aes(x = sq_norm)) +
  geom_histogram(aes(y = after_stat(density)),
    fill = "lightblue",
    bins = 100) +
  geom_density(col = "black",
    linewidth = 0.75) +
  labs(x = TeX("$||p||^2$"),
    y = "Frequency ",
    title = TeX("Distribution of Squared Euclidean Norm of $p$")) +
  theme_minimal() +
  theme(plot.title = element_text(hjust = 0.5),
    axis.title = element_text(color = "black"),
    axis.title.y = element_text(angle = 0, vjust = 0.5))
```

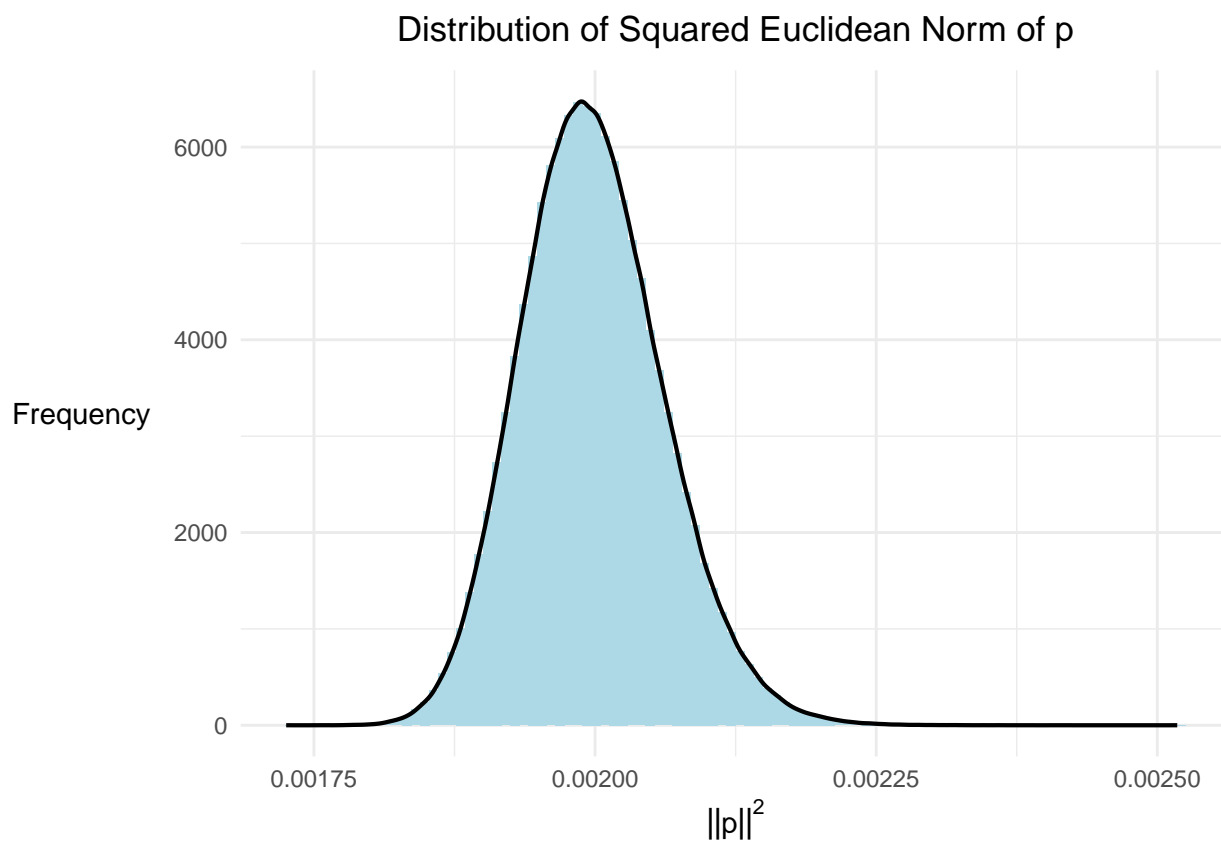


Comment.

- Now $n = 1000$:

```
probs <- prob_generator(N = 1000, m = 1000000)
sq_norms <- data.frame(sq_norm = colSums(probs^2))

sq_norms %>%
  ggplot(aes(x = sq_norm)) +
  geom_histogram(aes(y = after_stat(density)),
    fill = "lightblue",
    bins = 100) +
  geom_density(col = "black",
    linewidth = 0.75) +
  labs(x = TeX("$||p||^2$"),
    y = "Frequency ",
    title = TeX("Distribution of Squared Euclidean Norm of $p$")) +
  theme_minimal() +
  theme(plot.title = element_text(hjust = 0.5),
    axis.title = element_text(color = "black"),
    axis.title.y = element_text(angle = 0, vjust = 0.5))
```



Comment.