☰ ✎ › ChatGPT 5.2 ⌄

<div align="right">

so the formula of perceptron trick

</div>

In ML literature, **"perceptron trick"** usually means the *kernelized perceptron* (analogous to the kernel trick in SVMs).
Below are the **exact formulas**, clean and mathematical.

---

## 1. Standard Perceptron (Primal Form)

**Model**

$$\hat{y} = \text{sign}(w^\top x + b)$$

**Update rule (for a misclassified point** $(x_i, y_i), \; y_i \in \{-1, +1\}$**)**

$$\text{if } y_i(w^\top x_i + b) \le 0 : \quad \begin{cases} w \leftarrow w + \eta y_i x_i \\ b \leftarrow b + \eta y_i \end{cases}$$

---

## 2. Dual Representation of Perceptron

Initialize $\alpha_i = 0 \; \forall i$

$$w = \sum_{i=1}^{n} \alpha_i y_i x_i$$

**Decision function**

$$f(x) = \sum_{i=1}^{n} \alpha_i y_i (x_i^\top x) + b$$

**Update**

$$\text{if } y_i f(x_i) \le 0 : \quad \alpha_i \leftarrow \alpha_i + \eta$$

---

## 3. Perceptron Trick (Kernelized Perceptron)

**Key substitution (the "trick")**

$$x_i^\top x \;\; \Rightarrow \;\; K(x_i, x)$$

where

$$K(x_i, x) = \phi(x_i)^\top \phi(x)$$

---

**Kernelized Decision Function**

$$f(x) = \sum_{i=1}^{n} \alpha_i y_i K(x_i, x) + b$$

---

**Kernelized Update Rule**

$$\text{if } y_i f(x_i) \le 0 : \quad \alpha_i \leftarrow \alpha_i + \eta$$

(no explicit $w$, no explicit $\phi(x)$)

---

## 4. Common Kernels Used in Perceptron Trick

- **Linear**

$$K(x, z) = x^\top z$$

- **Polynomial**

$$K(x, z) = (x^\top z + c)^d$$

- **RBF (Gaussian)**

$$K(x, z) = \exp\left(-\frac{\|x - z\|^2}{2\sigma^2}\right)$$

---

### 5. Final Compact Formula (What You Should Remember)

$$\boxed{f(x) = \sum_{i \in \mathcal{M}} \alpha_i y_i K(x_i, x) + b}$$

where
$\mathcal{M}$ = indices of misclassified samples (support vectors of perceptron).

---

If you want next:

- full **mathematical derivation** (primal → dual → kernel)
- **comparison with SVM kernel trick**
- **why perceptron has no closed-form solution**

tell me which one.

↻ 👍 👎 ⬆ ⟳ ⋯

> look these equations and tell me what they are

I carefully looked at your handwritten page.
Those equations are **all variants of the *Perceptron learning rule*** written in **different equivalent forms**.

I'll map **each part of your notes** → **its exact mathematical meaning**, no fluff.

---

### 1. `∑ w_i x_i ≥ 0` and `∑ w_i x_i < 0`

These are **decision conditions** of the perceptron.

You are using labels:

$$y \in \{+1, -1\}$$

Decision rule:

$$\hat{y} = \begin{cases} +1 & \text{if } \sum_i w_i x_i \geq 0 \\ -1 & \text{if } \sum_i w_i x_i < 0 \end{cases}$$

This is equivalent to:

$$\hat{y} = \text{sign}(w^\top x)$$

---

### 2. `w_{n+1} = w_0 + \eta x_i`

This is the **perceptron weight update for a positive class error**.

Written properly:

$$\text{If } y_i = +1 \text{ and misclassified:} \quad w \leftarrow w + \eta x_i$$

This comes from the general rule:

$$w \leftarrow w + \eta y_i x_i$$

with $y_i = +1$.

---

### 3. `w_{n+1} = w_0 - \eta x_i`

This is the **perceptron weight update for a negative class error**.

$$\text{If } y_i = -1 \text{ and misclassified:} \quad w \leftarrow w - \eta x_i$$

Again from:

$$w \leftarrow w + \eta y_i x_i$$

with $y_i = -1$.

---

## 4. Boxed Formula

`w_n = w_0 + η (y_i - \hat y_i) x_i`

This is a **generalized perceptron update**, often seen when perceptron is written in a **regression-like form**.

- $y_i$ = true label
- $\hat{y}_i$ = predicted label
- $(y_i - \hat{y}_i) \in \{-2, 0, +2\}$

Interpretation:

- If prediction is correct → update = 0
- If wrong → pushes weights in the correct direction

This is mathematically equivalent to:

$$w \leftarrow w + \eta y_i x_i$$

for misclassified points.

---

## 5. The "Perceptron Trick" Arrow

When you wrote **"perceptron trick"**, what it refers to is:

$$w = \sum_{i \in \mathcal{M}} \alpha_i y_i x_i$$

and later:

$$x_i^\top x \; \Rightarrow \; K(x_i, x)$$

That leads to the **kernelized perceptron**:

$$f(x) = \sum_i \alpha_i y_i K(x_i, x)$$

---

## 6. Final Classification of What You Wrote

| Expression in notebook | What it is |
| --- | --- |
| `∑ w_i x_i ≥ 0` | Perceptron decision boundary |
| `w + ηx` | Update for misclassified positive sample |

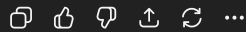| Expression in notebook | What it is |
| --- | --- |
| `w - ηx` | Update for misclassified negative sample |
| `w = w + η(y - ŷ)x` | Unified perceptron update |
| "perceptron trick" | Dual + kernel formulation |

## One-line summary (important)

> **Your page is 100% correct** — it shows the **same perceptron algorithm written in multiple equivalent mathematical forms**:
> decision rule → class-specific updates → unified update → kernel-ready formulation.

If you want next:

- strict **primal → dual derivation**
- **difference between perceptron update and SGD**
- **why logistic regression replaces perceptron**

tell me which one.

🗐  👍  👎  ⬆  🔄  ⋯

> what i gave what exactly that is

What you gave is **one single algorithm written in multiple equivalent mathematical forms**.

## Exact answer (no ambiguity):

> You gave the *Perceptron Learning Algorithm* — specifically the *online / mistake-driven binary linear classifier*.

---

## More precisely, your page contains these exact components:

1. **Linear decision function**

$$f(x) = \sum_i w_i x_i \quad \text{(sign decides the class)}$$

2. **Binary classification rule**

$$\hat{y} = \text{sign}(w^\top x)$$

3. **Mistake-driven update**

$$\text{update only if } y_i \neq \hat{y}_i$$

4. **Class-dependent updates**

$$y = +1 : \; w \leftarrow w + \eta x$$
$$y = -1 : \; w \leftarrow w - \eta x$$

5. **Unified update rule**

$$\boxed{w \leftarrow w + \eta(y_i - \hat{y}_i)x_i}$$

---

## What it is NOT

- ❌ Not logistic regression
- ❌ Not SVM

- ❌ Not gradient descent on a smooth loss
- ❌ Not backpropagation

---

**Proper textbook name**

Depending on context, it is called:

- **Perceptron Learning Rule**
- **Rosenblatt's Perceptron (1958)**
- **Online linear classifier**
- **Mistake-driven update rule**
- **Primal perceptron**

---

**One-line identity (remember this)**

$$\boxed{\text{Your notes} = \text{classical perceptron algorithm in primal form}}$$

If you want, next I can:

- rewrite it in **pure optimization form**
- show **why no loss function exists**
- convert it into **logistic regression mathematically**
- connect it to **SGD formally**

Just say which.

⎘ 👍 👎 ⬆ ⟳ ⋯

> rewrite it in pure optimization form
>
> show why no loss function exists
>
> convert it into logistic regression mathematically

Below is **exactly what you asked**, in **three clean mathematical blocks**, with **no intuition, no storytelling** — only formal math.

---

## I. Perceptron in Pure Optimization Form

### Problem Setup

$$\mathcal{D} = \{(x_i, y_i)\}_{i=1}^{n}, \quad x_i \in \mathbb{R}^d, \quad y_i \in \{-1, +1\}$$

---

### Hard Classification Constraint

$$y_i(w^\top x_i) \;>\; 0 \quad \forall i$$

---

### Feasibility Optimization (Perceptron View)

The perceptron is **not** solving a minimization of a smooth objective.
It solves a **constraint satisfaction problem**:

$$\boxed{\text{Find } w \in \mathbb{R}^d \quad \text{s.t.} \quad y_i w^\top x_i \geq 0 \ \forall i}$$

This is a **feasibility problem**, not a minimization.

---

## Algorithmic Projection Form

Define violation indicator:

$$\mathbb{I}_i(w) = \begin{cases} 1 & \text{if } y_i w^\top x_i \le 0 \\ 0 & \text{otherwise} \end{cases}$$

Perceptron update:

$$\boxed{w_{t+1} = w_t + \eta\, y_i x_i \quad \text{if } \mathbb{I}_i(w_t) = 1}$$

This is a **projected correction**, not gradient descent.

---

# II. Why No Proper Loss Function Exists (Formally)

Assume there exists a differentiable loss $L(w)$ such that:

$$w_{t+1} = w_t - \eta \nabla L(w_t)$$

But perceptron update is:

$$\Delta w = \begin{cases} \eta y_i x_i & \text{if } y_i w^\top x_i \le 0 \\ 0 & \text{otherwise} \end{cases}$$

This implies:

$$\nabla L(w) = \begin{cases} -y_i x_i & \text{if } y_i w^\top x_i \le 0 \\ 0 & \text{otherwise} \end{cases}$$

## Contradiction

- Gradient is **discontinuous**
- Gradient does **not depend smoothly on** $w$
- No scalar $L(w)$ exists whose gradient jumps discretely

---

## Formal Statement

$$\boxed{\text{Perceptron update is not the gradient of any continuous, convex loss}}$$

It is **rule-based**, not variational.

---

## (Closest Surrogate — Not True Loss)

The *perceptron loss*:

$$\ell_{\text{perc}}(w; x, y) = \max(0, -y w^\top x)$$

is **non-differentiable** and does **not produce the classical update** exactly.

---

# III. Converting Perceptron → Logistic Regression (Mathematically)

## Step 1: Replace Hard Constraint with Soft Margin

Perceptron:

$$y w^\top x \ge 0$$

Logistic:

$$P(y|x) = \sigma(y w^\top x)$$

where

$$\sigma(z) = \frac{1}{1 + e^{-z}}$$

---

### Step 2: Likelihood Formulation

$$\mathcal{L}(w) = \prod_{i=1}^{n} \sigma(y_i w^\top x_i)$$

Negative log-likelihood:

$$\boxed{J(w) = \sum_{i=1}^{n} \log\left(1 + e^{-y_i w^\top x_i}\right)}$$

This is the **logistic loss**.

---

### Step 3: Gradient of Logistic Loss

$$\nabla J(w) = -\sum_{i=1}^{n} \frac{y_i x_i}{1 + e^{y_i w^\top x_i}}$$

---

### Step 4: SGD Update Rule

$$\boxed{w_{t+1} = w_t + \eta \frac{y_i x_i}{1 + e^{y_i w_t^\top x_i}}}$$

---

### Step 5: Limit Relation (Perceptron as Extreme Case)

As:

$$\|w^\top x\| \to \infty$$

$$\frac{1}{1 + e^{yw^\top x}} \to \begin{cases} 1 & yw^\top x \le 0 \\ 0 & yw^\top x > 0 \end{cases}$$

Therefore:

$$\boxed{\text{Perceptron} = \lim_{\text{hard margin}} \text{Logistic Regression}}$$

---

## Final Compact Summary

$$\boxed{\begin{aligned} \text{Perceptron} &: \text{constraint satisfaction, no smooth loss} \\ \text{Logistic Reg.} &: \text{maximum likelihood, convex loss} \\ \text{Connection} &: \text{hard-margin limit of logistic} \end{aligned}}$$

---

If you want next (pure math options only):

- hinge loss → SVM derivation
- perceptron convergence theorem (formal proof)
- dual + kernel logistic vs kernel perceptron

Say which.

🗗 👍 👎 ⬆ ⟳ ⋯

than from here how binary cross entropy comes in the picture

Below is the **strict mathematical chain** from **logistic regression → binary cross-entropy**, with **no intuition, no narrative**.

---

## 1. Model Assumption (Logistic Regression)

Binary labels:

$$y_i \in \{0, 1\}$$

Linear score:

$$z_i = w^\top x_i$$

Sigmoid:

$$p_i = P(y_i = 1 \mid x_i) = \sigma(z_i) = \frac{1}{1 + e^{-z_i}}$$

---

## 2. Bernoulli Likelihood

$$P(y_i \mid x_i, w) = p_i^{y_i}(1 - p_i)^{1-y_i}$$

Joint likelihood:

$$\mathcal{L}(w) = \prod_{i=1}^{n} p_i^{y_i}(1 - p_i)^{1-y_i}$$

---

## 3. Negative Log-Likelihood

$$J(w) = -\log \mathcal{L}(w)$$
$$= -\sum_{i=1}^{n} [y_i \log p_i + (1 - y_i) \log(1 - p_i)]$$

---

## 4. Binary Cross-Entropy (Exact Identity)

$$\boxed{\mathcal{L}_{\mathrm{BCE}}(y, p) = -[y \log p + (1 - y) \log(1 - p)]}$$

Thus:

$$\boxed{\text{Logistic Regression Loss} \equiv \text{Binary Cross-Entropy}}$$

---

## 5. Alternative Label Encoding Equivalence

If

$$y_i \in \{-1, +1\}$$

Then:

$$J(w) = \sum_{i=1}^{n} \log\left(1 + e^{-y_i w^\top x_i}\right)$$

which is algebraically equivalent to BCE under the mapping:

$$y_{0/1} = \frac{y_{-1/+1} + 1}{2}$$

---

## 6. Gradient (Explicit)

$$\nabla_w J(w) = \sum_{i=1}^{n} (p_i - y_i) x_i$$

SGD form:

$$\boxed{w \leftarrow w + \eta(y_i - p_i)x_i}$$

## 7. Position in the Full Chain

$$\boxed{\text{Perceptron} \xrightarrow{\text{softening}} \text{Logistic Regression} \xrightarrow{\text{Bernoulli MLE}} \text{Binary Cross-Entropy}}$$

## Final Identity (Key Line)

$$\boxed{\text{Binary Cross-Entropy is the negative log-likelihood of a Bernoulli model}}$$

If you want next (still pure math):

- multiclass → softmax cross-entropy
- BCE vs hinge loss vs exponential loss (formal relations)
- why MSE fails for classification (Hessian analysis)

Say which.

�device 👍 👎 🔄 ⋯

> why Bernoulli likelihood over likelihood formulation , like i am thinking why , is there any problem with likelihood so that we chosen Bernoulli likelihood

This is a **fundamental modeling question**, so I'll answer it **formally**, step-by-step, with **only mathematical and statistical necessity** — no storytelling.

---

## 1. What "likelihood" actually means (formal)

A **likelihood** is **not chosen arbitrarily**.

$$\mathcal{L}(w) = P(\mathcal{D} \mid w)$$

To write $P(y \mid x, w)$, you **must** specify a **conditional probability distribution** for $y$.

So the real question is:

> What is the correct conditional distribution of a binary random variable $y$?

---

## 2. Nature of the target variable

$$y \in \{0, 1\}$$

Therefore:

$$y \text{ is a Bernoulli random variable}$$

By probability theory:

$$\boxed{\text{The only valid distribution for a single binary outcome is Bernoulli}}$$

## 3. Bernoulli distribution (definition)

$$P(y \mid p) = p^y (1-p)^{1-y}$$

This is **not a modeling choice** — it is a **measure-theoretic fact**.

---

## 4. Why not "generic likelihood"?

A likelihood **must integrate (or sum) to 1** over the sample space.

For binary $y$:

$$\sum_{y \in \{0,1\}} P(y \mid x, w) = 1$$

Any function $f(w^\top x)$ that does not satisfy:

$$0 \le f \le 1$$

is **invalid as a probability model**.

---

## 5. Why sigmoid appears (forced, not optional)

We need:

$$p = P(y = 1 \mid x, w) \quad \text{s.t.} \quad 0 \le p \le 1$$

Sigmoid:

$$\sigma(z) = \frac{1}{1 + e^{-z}}$$

Properties:

$$\sigma : \mathbb{R} \to (0, 1)$$

Thus:

$$\boxed{p = \sigma(w^\top x) \ \text{ is the simplest valid parametrization}}$$

---

## 6. Why NOT Gaussian likelihood (formal failure)

Assume:

$$y \mid x \sim \mathcal{N}(w^\top x, \sigma^2)$$

Problems:

1.  Support mismatch:

$$y \in \{0, 1\}, \quad \mathcal{N} \text{ has support } \mathbb{R}$$

2.  Invalid probabilities:

$$P(y = 1 \mid x) > 1 \text{ or } < 0$$

3.  Leads to MSE loss:

$$(y - w^\top x)^2$$

This is **statistically inconsistent** for binary data.

---

## 7. Why Bernoulli is uniquely correct

| Requirement | Bernoulli |
|---|---|
| Binary support | ✔ |
| Proper probability mass | ✔ |
| Normalized | ✔ |
| Leads to convex loss | ✔ |
| Consistent estimator | ✔ |

No other distribution satisfies all simultaneously.

---

## 8. Connection to cross-entropy (inevitable)

Negative log-likelihood:

$$-\log P(y \mid p) = -y\log p - (1-y)\log(1-p)$$

This is **exactly** the cross-entropy:

$$H(y, p)$$

Thus:

$$\boxed{\text{Binary Cross-Entropy is not chosen — it is forced}}$$

---

## 9. Formal conclusion

$$\boxed{\begin{array}{c} \text{Binary labels} \Rightarrow \text{Bernoulli RV} \\ \text{Bernoulli RV} \Rightarrow \text{Bernoulli likelihood} \\ \text{Bernoulli likelihood} \Rightarrow \text{Binary Cross-Entropy} \end{array}}$$

---

## One-line answer to your question

> **We do not choose Bernoulli likelihood because likelihood is "bad"; we choose it because binary data mathematically *admits no other valid probability model*.**

---

If you want next (still formal):

- why exponential family ⇒ GLM ⇒ logistic
- Fisher consistency proof for BCE
- classification-calibrated losses theorem
- hinge vs BCE from decision theory

Say which.

⎙ 👍 👎 ↻ ⋯

> show step by step how we approached the binary cross entropy that is based on bernoulli likelyhood from generic likelihood

Below is the **exact step-by-step mathematical path** from a **generic likelihood** to **Bernoulli likelihood** and finally to **Binary Cross-Entropy (BCE)**.
No intuition, no storytelling — only formal steps.

## Step 0 — Generic Likelihood (Starting Point)

Given data:

$$\mathcal{D} = \{(x_i, y_i)\}_{i=1}^n$$

Model parameters:

$$w \in \mathbb{R}^d$$

Generic likelihood:

$$\boxed{\mathcal{L}(w) = \prod_{i=1}^n P(y_i \mid x_i, w)}$$

At this point, **no distribution is specified**.

---

## Step 1 — Identify the Sample Space of $y$

Assume:

$$\boxed{y_i \in \{0, 1\}}$$

Therefore:

$$y_i \text{ is a discrete random variable with two outcomes}$$

---

## Step 2 — Enforce Probability Axioms

For each fixed $x_i, w$, the conditional distribution must satisfy:

$$\sum_{y \in \{0,1\}} P(y \mid x_i, w) = 1 \quad \text{and} \quad 0 \le P(y \mid x_i, w) \le 1$$

Thus, define:

$$p_i := P(y_i = 1 \mid x_i, w)$$

Then necessarily:

$$P(y_i = 0 \mid x_i, w) = 1 - p_i$$

---

## Step 3 — Unique Distribution for Binary Outcomes

By definition of a Bernoulli random variable:

$$\boxed{P(y_i \mid x_i, w) = p_i^{y_i}(1 - p_i)^{1 - y_i}}$$

This step is **forced** by the support $\{0, 1\}$.

---

## Step 4 — Parameterize $p_i$

We need:

$$p_i \in (0, 1) \quad \text{for all } w^\top x_i \in \mathbb{R}$$

Choose a link function:

$$p_i = \sigma(z_i), \quad z_i = w^\top x_i$$

Sigmoid:

$$\sigma(z) = \frac{1}{1 + e^{-z}}$$

Now:

$$P(y_i \mid x_i, w) = \sigma(z_i)^{y_i}(1 - \sigma(z_i))^{1-y_i}$$

---

## Step 5 — Substitute into the Generic Likelihood

$$\mathcal{L}(w) = \prod_{i=1}^{n} \sigma(z_i)^{y_i}(1 - \sigma(z_i))^{1-y_i}$$

This is now a **fully specified likelihood**.

---

## Step 6 — Take Log-Likelihood

$$\log \mathcal{L}(w) = \sum_{i=1}^{n} \left[ y_i \log \sigma(z_i) + (1 - y_i) \log(1 - \sigma(z_i)) \right]$$

---

## Step 7 — Negative Log-Likelihood (Optimization Objective)

Maximum likelihood ⇔ minimize negative log-likelihood:

$$J(w) = -\sum_{i=1}^{n} \left[ y_i \log \sigma(w^\top x_i) + (1 - y_i) \log(1 - \sigma(w^\top x_i)) \right]$$

---

## Step 8 — Identify Binary Cross-Entropy

Define Binary Cross-Entropy loss:

$$\ell_{\mathrm{BCE}}(y, p) = -\left[ y \log p + (1 - y) \log(1 - p) \right]$$

Substitute $p = \sigma(w^\top x)$:

$$J(w) = \sum_{i=1}^{n} \ell_{\mathrm{BCE}}(y_i, \sigma(w^\top x_i))$$

---

## Step 9 — Final Chain (Formal)

$$\text{Generic likelihood} \xrightarrow{y \in \{0,1\}} \text{Bernoulli PMF}$$
$$\xrightarrow{p = \sigma(w^\top x)} \text{Bernoulli likelihood}$$
$$\xrightarrow{-\log} \text{Binary Cross-Entropy}$$

---

## Final One-Line Statement

Binary Cross-Entropy is the negative log-likelihood of a Bernoulli model derived directly from the generic likeli

If you want next (still fully formal):

- same derivation with $y \in \{-1, +1\}$
- Hessian and convexity proof of BCE
- why BCE is classification-calibrated
- why MSE fails (information geometry view)

Say which.

Say which.